

Bayesian Nonparametric Ordination for the Analysis of Microbial Communities

Boyu Ren¹, Sergio Bacallado², Stefano Favaro³, Susan Holmes⁴ and
Lorenzo Trippa¹

¹Harvard University, Cambridge, US

²University of Cambridge, Cambridge, UK

³Università degli Studi di Torino and Collegio Carlo Alberto, Turin,
Italy

⁴Stanford University, Stanford, US

January 21, 2016

Abstract

Human microbiome studies use sequencing technologies to measure the abundance of bacterial species or Operational Taxonomic Units (OTUs) in samples of biological material. Typically the data are organized in contingency tables with OTU counts across heterogeneous biological samples. In the microbial ecology community, ordination methods are frequently used to investigate latent factors or clusters that capture and describe variations of OTU counts across biological samples. It remains important to evaluate how uncertainty in estimates of each biological sample's microbial distribution propagates to ordination analyses, including visualization of clusters and projections of biological samples on low dimensional spaces. We propose a Bayesian analysis for dependent distributions to endow frequently used ordinations with estimates of uncertainty. A Bayesian nonparametric prior for dependent normalized random measures is constructed, which is marginally equivalent to the normalized generalized Gamma process, a well-known prior for nonparametric analyses. In our prior the dependence and similarity between microbial distributions is represented by latent factors that concentrate in a low dimensional space. We use a shrinkage prior to tune the dimensionality of the latent factors. The resulting posterior samples of model parameters can be used to evaluate uncertainty in analyses routinely applied in microbiome studies. Specifically, by combining them with multivariate data analysis techniques we can visualize credible regions in ecological ordination plots. The characteristics of the proposed model are illustrated through a simulation study and applications in two microbiome datasets.

1 Introduction

Next generation sequencing (NGS) has transformed the study of microbial ecology. Through the availability of cheap efficient amplification and sequencing, marker genes such as 16S rRNA are used to provide inventories of bacteria in many different environments. For instance soil and waste water microbiota have been inventoried (DeSantis et al., 2006) as well as the human body (Dethlefsen et al., 2007). NGS also enables researchers to describe the *metagenome* by computing counts of DNA reads and matching them to the genes present in various environments.

Over the last ten years, numerous studies have shown the effects of environmental and clinical factors on the bacterial communities of the human microbiome. These studies enhance our understanding of how the microbiome is involved in obesity (Turnbaugh et al., 2009), Crohn’s disease (Quince et al., 2013) or diabetes (Kostic et al., 2015). Studies are currently underway to improve our understanding of the effects of antibiotics (Dethlefsen and Relman, 2011), pregnancy (DiGiulio et al., 2015) and other perturbations to the human microbiome.

Common microbial ecology pipelines either start by grouping the 16S rRNA sequences into known Operational Taxonomic Units (OTUs) or taxa as done in Caporaso et al. (2010) or denoising and grouping the reads into more refined strains sometimes referred to as oligotypes or phylotypes (Rosen et al., 2012; Eren et al., 2014). We will call both types of groupings OTUs to maintain consistency. In all cases the data are analyzed in the form of contingency tables of read counts per sample for the different OTUs. Associated to these contingency tables are clinical and environmental covariates such as time, treatment and patients’ BMI, information collected on the same biological samples or environments. These are sometimes misnamed “metadata”; this contiguous information is usually fundamental in the analyses. The data are often assembled in multi-type structures, for instance `phyloseq` (McMurdie and Holmes, 2013) uses lists (S4 classes) to capture all the different aspects of the data at once.

Currently bioinformaticians and statisticians analyze the preprocessed microbiome data using linear ordination methods such as Correspondence Analysis (CA), Canonical or Constrained Correspondence Analysis (CCA) and Multidimensional Scaling (MDS) (Caporaso et al., 2010; Oksanen et al., 2015; McMurdie and Holmes, 2013). Distance-based ordination methods use measures of between-sample or Beta diversity, such as the Unifrac distance (Lozupone and Knight, 2005). These analyses can reveal clustering of biological samples or taxa, or meaningful ecological or clinical gradients in the community structure of the bacteria. Clustering, when it occurs indicates a latent variable which is discrete, whereas gradients correspond to latent continuous variables. Following these exploratory stages, confirmatory analyses can include differential abundance testing (McMurdie and Holmes, 2014), two-sample tests for Beta diversity scores (Anderson et al., 2006), ANOVA permutation tests in CCA (Oksanen et al., 2015), or tests based on generalized linear models that include adjustment for multiple confounders (Paulson et al., 2013).

The interaction between these tasks can be problematic. In particular, the uncertainty in the estimation of OTUs’ prevalence is often not propagated to subsequent

steps (Peiffer et al., 2013). Moreover, unequal sequencing depths generate variations of the number of OTUs with zero counts across biological samples. Finally, the hypotheses tested in the inferential step are often formulated after significant exploration of the data, and are sensitive to earlier choices in data preprocessing.

These issues motivate a Bayesian approach that enables us to integrate the steps of the analytical pipeline. Holmes et al. (2012) have suggested the use of a simple Dirichlet-Multinomial model for these data; however, in that analysis the multinomial probabilities for each sample are independent in the prior and posterior, which fails to capture underlying relationships between biological samples. On the other hand, the correlation between OTUs is negative in both prior and posterior, which is not consistent with the fact that some OTUs can only exist in tight-knit communities.

We propose a Bayesian procedure, which jointly models the read counts from different OTUs and sample-specific latent multinomial distributions, allowing for correlations between OTUs. The prior assigned to these multinomial probabilities is highly flexible, such that the analysis learns the dependence structure from the data, rather than constraining it *a priori*. The method can deal with uncertainty coherently, provides model-based visualizations of the data and describes the effects of environmental covariates.

Bayesian analysis with Dirichlet priors is a suitable starting point for microbiome data, since the OTUs distributions are inherently discrete. Moreover, Bayesian non-parametric priors for discrete distributions, suitable for an unbounded number of OTUs, have been the topic of intense research in recent years. General classes of priors such as normalized random measures have been developed, and their properties in relation to classical estimators of species diversity are well-understood (Ferguson, 1973; Lijoi and Prünster, 2010). The problem of modeling dependent distributions has also been extensively studied since the proposal of the Dependent Dirichlet Process (MacEachern, 2000) by Müller et al. (2004); Rodríguez et al. (2009) and Griffin et al. (2013).

In this paper, we try to capture the variation in the composition of microbial communities as a result of the observed and unobserved samples’ characteristics. With this goal we introduce a model which expresses the dependence between OTUs abundances in different environments through vectors embedded in a low dimensional space. Our model has aspects in common with nonparametric priors for dependent distributions, including a generalized Dirichlet type marginal prior on each distribution, but is also similar in spirit to the multivariate methods currently employed in the microbial ecology community. Namely, it allows us to visualize the relationship between biological samples through low dimensional projections.

The paper is organized as follows, Section 2 describes a prior for dependent microbial distributions, first constructing the marginal prior of a single distribution through manipulation of a Gaussian process, and then extending this to multiple correlated distributions. The extension is achieved through a set of continuous latent factors, one for each biological sample, whose prior is one frequently used in Bayesian factor analyses. Section 3 derives an MCMC sampling algorithm for posterior inference and a fast algorithm to estimate biological samples’ similarity. Section 4 discusses a method

for visualizing the uncertainty in ordinations through conjoint analysis. Section 5 contains analyses of simulated data, which serve to demonstrate desirable properties of the method, followed by applications to real microbiome data in Section 6. Section 7 discusses potential improvement and concludes. The code for implementing the analyses discussed in this article can be found in the repository `DirichletFactor`.

2 Probability Model

We construct discrete distributions $\{P^j; j \in \mathcal{J}\}$ indexed by biological samples in a set \mathcal{J} . The distributions are supported on a common, countable set of OTUs, Z_1, Z_2, \dots , in a measure space $(\mathcal{Z}, \mathcal{F})$, where \mathcal{Z} is assumed to be complete and separable. Every OTU Z_i is associated to parameters $\sigma_i \in (0, 1)$ and $\mathbf{X}_i \in \mathbb{R}^\infty$. Every microbial distribution P^j is associated to a parameter $\mathbf{Y}^j \in \mathbb{R}^\infty$. The measure P^j is defined by

$$\begin{aligned} P^j(A) &= M^j(A)/M^j(\mathcal{Z}), \\ M^j(A) &= \sum_{i=1}^{\infty} I(Z_i \in A) \sigma_i \langle \mathbf{X}_i, \mathbf{Y}^j \rangle^{+2}, \end{aligned} \tag{1}$$

for every set $A \in \mathcal{F}$. Here $I(\cdot)$ and $\langle \cdot, \cdot \rangle$ are the standard indicator and inner product functions, while x^+ is the positive part of x .

In subsection 2.1 we consider a single microbial distribution P^j with fixed parameter \mathbf{Y}^j , and define a prior on $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \dots)$ and $(\mathbf{X}_i)_{i \geq 1}$ which makes P^j a Dirichlet process (Ferguson, 1973). The degree of similarity between the discrete distributions $\{P^j; j \in \mathcal{J}\}$ is summarized by the Gram matrix $\{\phi(j, j') = \langle \mathbf{Y}^j, \mathbf{Y}^{j'} \rangle; j, j' \in \mathcal{J}\}$. Subsection 2.2 discusses the interpretation of this matrix, and subsection 2.3 proposes a prior for the parameters $\{\mathbf{Y}^j, j \in \mathcal{J}\}$ which has been previously used in Bayesian factor analysis, and which has the effect of shrinking the dimensionality of the Gram matrix ϕ . The parameters $\{\mathbf{Y}^j, j \in \mathcal{J}\}$ or ϕ can be used to visualize and understand variations of microbial distributions across biological samples.

2.1 Construction of a Dirichlet Process

The prior on $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \dots)$ is the distribution of ordered points $(\sigma_i > \sigma_{i+1})$ in a Poisson process on $(0, 1)$ with intensity

$$\nu(\sigma) = \alpha \sigma^{-1} (1 - \sigma)^{-1/2}, \tag{2}$$

where $\alpha > 0$ is a concentration parameter. Fix j , and let $\mathbf{Y}^j = (Y_{l,j}, l \geq 1)$ be a fixed sequence of real numbers with $\sum_l Y_{l,j}^2 = 1$. Using the notation $\mathbf{X}_i = (X_{1,i}, X_{2,i}, \dots)$, we let $X_{l,i}$, $i, l = 1, 2, \dots$ be independent and $N(0, 1)$ *a priori*.

Finally, let G be a nonatomic probability measure on $(\mathcal{Z}, \mathcal{F})$, and Z_1, Z_2, \dots be a sequence of independent random variables with distribution G . We claim that the probability distribution P^j defined in Equation (1) is a Dirichlet Process with base measure G .

We note that the point process σ defines an infinite sequence of positive numbers, the products $\langle \mathbf{X}_i, \mathbf{Y}^j \rangle$, $i = 1, 2, \dots$, are independent Gaussian $N(0, 1)$ variables, and that the intensity ν satisfies the inequality $\int_0^1 \sigma d\nu < \infty$. These facts directly imply that with probability 1, $0 < M^j(A) < \infty$ when $G(A) > 0$. It also follows that for any sequence of disjoint sets $A_1, A_2, \dots \in \mathcal{F}$ the corresponding random variables $M^j(A_i)$'s are independent. In different words, M^j is a completely random measure (Kingman, 1967). The marginal Lévy intensity can be factorized as $\mu_M(ds) \times G(dz)$, where

$$\begin{aligned} \mu_M(ds) &\propto \int_0^1 \nu(\sigma) \left(\frac{1}{\sigma} \right)^{1/2} s^{-1/2} \exp\left(-\frac{s}{2\sigma}\right) d\sigma ds \\ &\propto \frac{\exp(-s/2)}{s} ds, \quad \text{for } s \in (0, \infty). \end{aligned}$$

The above expression shows that M^j is a Gamma process. We recall that the Lévy intensity of a Gamma process is proportional to the map $s \mapsto \exp(-c \times s) \times s^{-1}$, where c is a positive scale parameter. In Ferguson (1973) it is shown that a Dirichlet process can be defined by normalizing a Gamma process. It directly follows that P^j is a Dirichlet Process with base measure G .

Remark. *Our construction can be extended to a wider class of normalized random measures (James, 2002; Regazzini et al., 2003) by changing the intensity ν that defines the Poisson process σ . If we set*

$$\nu(\sigma) = \alpha \sigma^{-1-\beta} (1 - \sigma)^{-1/2+\beta},$$

$\beta \in [0, 1)$, in our definition of M^j , then the Lévy intensity of the random measure in (1) becomes proportional to

$$s^{-1-\beta} \exp(-s/2).$$

In this case the Lévy intensity indicates that M^j is a generalized Gamma process (Brix, 1999). We recall that by normalizing this class one obtains normalized generalized Gamma processes (Lijoi et al., 2007), which include the Dirichlet process and the normalized Inverse Gaussian process (Lijoi et al., 2005) as special cases.

A few comments capture the relation between our definition of $P^j(A)$ in (1) and alternative definitions of the Dirichlet Process. If we normalize m independent $\text{Gamma}(\alpha/m, 1/2)$ variables, we obtain a vector with $\text{Dirichlet}(\alpha/m, \dots, \alpha/m)$ distribution. To interpret our construction we can note that, when $\alpha/m < 1/2$, each of the $\text{Gamma}(\alpha/m, 1/2)$ components can be obtained by multiplying a $\text{Beta}(\alpha/m, 1/2 - \alpha/m)$ variable and an independent $\text{Gamma}(1/2, 1/2)$. The distribution of the $\langle \mathbf{X}_i, \mathbf{Y}^j \rangle^{+2}$ variables in (1) is in fact a mixture with a $\text{Gamma}(1/2, 1/2)$ component and a point mass at zero. Finally if we let m increase to ∞ , the law of the ordered $\text{Beta}(\alpha/m, 1/2 - \alpha/m)$ converges weakly to the law of ordered points of a Poisson point process on $(0, 1)$ with intensity ν (see Appendix A).

2.2 Dependent Dirichlet Processes

We use the representation for Dirichlet processes from Equation (1) to define a family of dependent Dirichlet processes. Let \mathcal{J} be a complete and separable index space and $\phi : \mathcal{J} \times \mathcal{J} \rightarrow (-1, 1)$ a strictly positive-definite kernel with $\phi(j, j) = 1$ for $j \in \mathcal{J}$. The index space can be for example the space of integer numbers \mathbb{N}^+ . By Mercer's theorem (Mercer, 1909), there exists a set of sequences $\{\mathbf{Y}^j = (Y_{1,j}, Y_{2,j}, \dots)\}_{j \in \mathcal{J}}$, such that $\sum_l Y_{l,j} Y_{l,j'} = \phi(j, j')$ for every pair $j, j' \in \mathcal{J}$. Geometrically $\phi(j, j')$ is the cosine of the angle between \mathbf{Y}^j and $\mathbf{Y}^{j'}$. A family of dependent Dirichlet processes is defined by setting

$$P^j(A) = \frac{\sum_i I(Z_i \in A) \times \sigma_i \langle \mathbf{X}_i, \mathbf{Y}^j \rangle^{+2}}{\sum_i \sigma_i \langle \mathbf{X}_i, \mathbf{Y}^j \rangle^{+2}}, \quad \forall j \in \mathcal{J}, \quad (3)$$

for every $A \in \mathcal{F}$. Here the sequence (Z_1, Z_2, \dots) and the array $(\mathbf{X}_1, \mathbf{X}_2, \dots)$, as in Section 2.1, contain independent and identically distributed random variables, while σ is our Poisson process on the unit interval. We will use the notation $Q_{i,j} = \langle \mathbf{X}_i, \mathbf{Y}^j \rangle$. This construction has an interpretable dependency structure between the P^j 's that we state in the next proposition.

Proposition 1. *There exists a real function $\eta : [0, 1] \rightarrow [0, 1]$ such that the correlation between $P^j(A)$ and $P^{j'}(A)$ is equal to $\eta(\phi(j, j'))$ for every A that satisfies $G(A) > 0$. In different words, the correlation between $P^j(A)$ and $P^{j'}(A)$ does not depend on the specific measurable set A , it is a function of the angle defined by \mathbf{Y}^j and $\mathbf{Y}^{j'}$.*

The proof is in Appendix B. The degree of similarity between random probability measures in the definition (3) is specified through the kernel ϕ . The first panel of Figure 1 shows a simulation of P^j 's. In this figure $\mathcal{J} = \{1, 2, 3, 4\}$. When ϕ , the cosine of the angle between two vectors \mathbf{Y}^j and $\mathbf{Y}^{j'}$, corresponding to distinct biological samples j and j' , decreases to -1 the random measures tend to concentrate on two disjoint sets. The second panel shows the function η that maps the $\phi(j, j')$'s into the correlations $\text{corr}(P^j(A), P^{j'}(A)) = \eta(\phi(j, j'))$. As expected the correlation increases with $\phi(j, j')$.

The next proposition provides mild conditions that guarantee a large support for the dependent Dirichlet processes that we defined. We assume here $\mathcal{Z} \subset \mathbb{R}$. This will be sufficient for most applications. The proof is in Appendix C.

Proposition 2. *Consider a collection of probability measures $(F_i, i = 1, \dots, d)$ on \mathcal{Z} , $\mathcal{J} = \{1, \dots, d\}$, a positive definite kernel ϕ and assume that the support of G coincides with \mathcal{Z} . The prior distribution in (3) assigns strictly positive probability to the neighborhood $\{(F'_1, \dots, F'_d) : |\int f_i dF'_j - \int f_i dF_j| < \epsilon, i = 1, \dots, m, j = 1, \dots, d\}$, where $\epsilon > 0$ and $f_i, i = 1, \dots, d$, are bounded continuous functions.*

In what follows we will replace the constraint $\sum_l Y_{l,j}^2 = 1$ with the requirement $\sum_l Y_{l,j}^2 < \infty$. The two constraints are equivalent for our purpose, because we normalize $M^j(\cdot) = \sum_i I(Z_i \in \cdot) \times \sigma_i \langle \mathbf{X}_i, \mathbf{Y}^j \rangle^{+2}$, and $\sum_l Y_{l,j}^2$ can be viewed as a scale parameter.

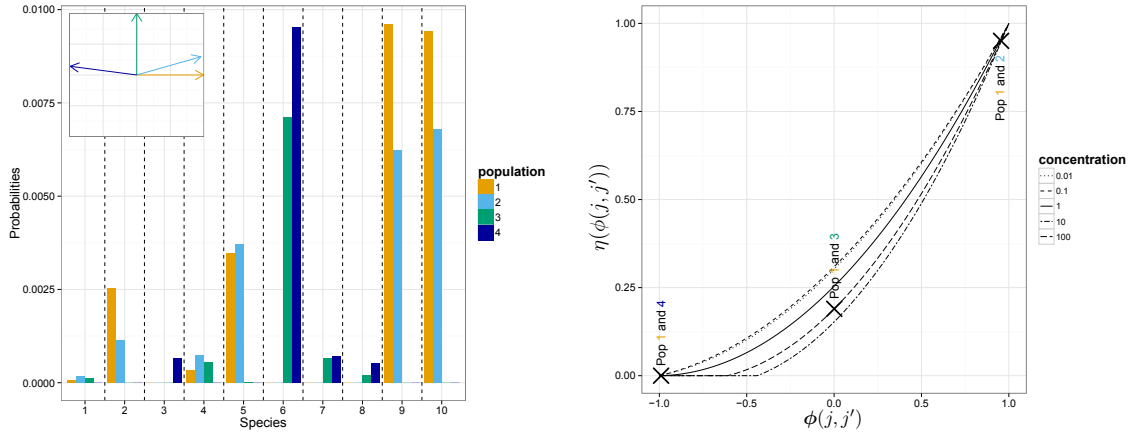


Figure 1: **Left panel:** realization of 4 microbial distributions from our dependent Dirichlet processes. We illustrate 10 representative OTUs and set $\alpha = 100$. The miniature figure at the top-left corner shows the relative positions of the four microbial distributions' vectors \mathbf{Y} . The OTUs are those associated to the 10 largest σ 's. As suggested by this panel, the larger the angle between \mathbf{Y} directions of two biological samples, the more the corresponding random distributions tend to concentrate on distinct sets. **Right panel:** correlation of two random probability measures when the cosine $\phi(j, j')$ between \mathbf{Y}^j and $\mathbf{Y}^{j'}$ varies from -1 to 1 . We consider five different values of the concentration parameter α . In the right panel we also mark with crosses the correlations between $P^j(A)$ and $P^{j'}(A)$ for pairs of biological samples j, j' considered in the left panel.

2.3 Prior on biological sample parameters

This subsection deals with the task of estimating the parameters $\mathbf{Y}^j, j \in \mathcal{J} = \{1, \dots, J\}$, that capture most of the variability observed when comparing J biological samples with different OTU counts. We define a joint prior on these factors which makes them concentrate on a low dimensional space; equivalently, the prior tends to shrink the nuclear norm of the Gram matrix $(\phi(j_1, j_2))_{j_1, j_2 \in \mathcal{J}}$. The problem of estimating low dimensional factor loadings or a low-rank covariance matrix is common in Bayesian factor analysis, and the prior defined below has been used in this area of research.

The parameters \mathbf{Y}^j can be interpreted as key characteristics of the biological samples that affect the relative abundance of OTUs. As in factor analysis, it is difficult to interpret these parameters unambiguously (Press and Shigemitsu, 1989; Rowe, 2002); however, the angles between their directions have a clear interpretation. As observed in Figure 1, if the kernel $\phi(j_1, j_2) \approx \sqrt{\phi(j_1, j_1)\phi(j_2, j_2)}$, the two microbial distributions P^{j_1} and P^{j_2} will be very similar. If $\phi(j_1, j_2) \approx 0$, then there will be little correlation between OTUs' abundances in the two samples. If $\phi(j_1, j_2) \approx -\sqrt{\phi(j_1, j_1)\phi(j_2, j_2)}$, then the two microbial distributions are concentrated on disjoint sets. This interpretation suggests PCA of the Gram matrix $(\phi(j_1, j_2))_{j_1, j_2 \in \mathcal{J}}$ as a useful exploratory data analysis technique.

It is common in factor analysis to restrict the dimensionality of factor loadings. In our model, this can be accomplished by choosing a number of degrees of freedom $m \leq J$, setting $(Y_{m+1,j}, Y_{m+2,j}, \dots)$ to be zero and adding an error term ϵ in the definition of $Q_{i,j}$, the OTU-specific latent weights,

$$Q_{i,j} = \langle \mathbf{X}_i, \mathbf{Y}^j \rangle + \epsilon_{i,j}, \quad (4)$$

where the $\epsilon_{i,j}$ are independent Normal variables. Recall that each sample-specific random distribution P^j is obtained by normalizing the random variables $\sigma_i(Q_{i,j}^+)^2$.

In most applications the dimensionality m is unknown. Several approaches to estimate m have been proposed (Lopes and West, 2004; Lee and Song, 2002; Lucas et al., 2006; Carvalho et al., 2008; Ando, 2009). However, most of them involve either calculation of Bayes Factors or complex MCMC algorithms. Instead we use a Normal shrinkage prior proposed by Bhattacharya and Dunson (2011). This prior includes an infinite sequence of factors ($m = \infty$), but the variability captured by this sequence of latent factors rapidly decreases to zero. A key advantage of the model is that it does not require to choose the number of factors. The prior is designed to replace direct selection of m with the strictly related goal of shrinking toward zero the unnecessary latent factors. In addition, this prior is nearly conjugate, which simplifies computations. The prior is defined as follows,

$$\begin{aligned} \gamma_l &\sim \text{Gamma}(a_l, 1), & \gamma'_{l,j} &\sim \text{Gamma}(v/2, v/2), \\ Y_{l,j} | \gamma_l, \gamma'_{l,j} &\sim N \left(0, (\gamma'_{l,j})^{-1} \prod_{k \leq l} \gamma_k^{-1} \right), & l \geq 1, j \in \mathcal{J}, \end{aligned} \quad (5)$$

where the random variables $\gamma = (\gamma_l, \gamma'_{l,j}; l, j \geq 1)$ are independent and, conditionally on these variables, the $Y_{l,j}$'s are independent.

When $a_l > 1$, the shrinkage strength *a priori* increases with the index l , and therefore the variability captured by each latent factor tends to decrease with l . We refer to Bhattacharya and Dunson (2011) for a detailed analysis of the prior in (5). In practice, the assumption of infinitely many factors is replaced for data analysis and posterior computations by a finite and sufficiently large number m of factors. This prior model is conditionally conjugate when paired with the dependent Dirichlet processes prior in subsection 2.2, a relevant and convenient characteristic for posterior simulations.

3 Posterior Analysis

Given an exchangeable sequence W_1, \dots, W_n from $P^j = M^j \times M^j(\mathcal{Z})^{-1}$ as defined in subsection 2.1, we can rewrite the likelihood function using variable augmentation as in James et al. (2009),

$$\prod_{i=1}^n P^j(\{W_i\}) = \int_0^\infty \frac{\exp[-M^j(\mathcal{Z}) T] \times T^{n-1}}{\Gamma(n)} \prod_{i=1}^I M^j(\{W_i^*\})^{n_i} dT. \quad (6)$$

Here W_1^*, \dots, W_I^* is the list of distinct values in (W_1, \dots, W_n) and n_1, \dots, n_I are the occurrences in (W_1, \dots, W_n) , so that $\sum_{i=1}^I n_i = n$. We use expression (6) to specify an algorithm that allows us to infer microbial abundances P^1, \dots, P^J in J biological samples.

We proceed, similarly to Muliere and Tardella (1998) and Ishwaran and James (2001), using truncated versions of the processes in subsection 2.2. We replace $\sigma = \{\sigma_i, i \geq 1\}$ with a finite number I of independent $\text{Beta}(\epsilon_I, 1/2 - \epsilon_I)$ points in $(0, 1)$. Appendix A shows that when I diverges, and $\epsilon_I = \alpha/I$, this finite dimensional version converges weakly to the process in (2). Each point σ_i is paired with a multivariate normal $\mathbf{Q}_i = (Q_{i,1}, \dots, Q_{i,J})$ with mean zero and covariance Σ . The distribution of $M_{i,j} = \sigma_i(Q_{i,j}^+)^2$ is a mixture of a point mass at zero and a Gamma distribution. In this section \mathbf{Q} and σ are finite dimensional, and the normalized vectors P^j , which assign random probabilities to I OTUs in J biological samples, are proportional to $(M_{1,j}, \dots, M_{I,j})$, $j = 1, \dots, J$. Note that P^j conditional on $1(Q_{1,j} > 0), \dots, 1(Q_{I,j} > 0)$ follows a Dirichlet distribution with parameters $\epsilon_I \times 1(Q_{1,j} > 0), \dots, \epsilon_I \times 1(Q_{I,j} > 0)$.

The algorithm is based on iterative sampling from the full conditional distributions. We first provide a description assuming that Σ is known. We then extend the description to allow sampling under the shrinkage prior in Section 2.3 and to infer Σ .

With I OTUs and J biological samples, the typical dataset is $\mathbf{n} = (\mathbf{n}_1, \dots, \mathbf{n}_J)$, where $\mathbf{n}_j = (n_{1,j}, \dots, n_{I,j})$ and $n_{i,j}$ is the absolute frequency of the i^{th} OTU in the j^{th} biological sample. We use the notation $n^j = \sum_{i=1}^I n_{i,j}$, $n_i = \sum_{j=1}^J n_{i,j}$, $\sigma = (\sigma_1, \dots, \sigma_I)$, $\mathbf{Y} = (\mathbf{Y}^j, j = 1, \dots, J)$ and $\mathbf{Q} = (Q_{i,j}, 1 \leq i \leq I, 1 \leq j \leq J)$. By using

the representation in (6) we introduce the latent random variables $\mathbf{T} = (T_1, \dots, T_J)$ and rewrite the posterior distribution of $(\boldsymbol{\sigma}, \mathbf{Q})$:

$$p(\boldsymbol{\sigma}, \mathbf{Q}|\mathbf{n}) \propto \left(\prod_{j=1}^J \prod_{i=1}^I (\sigma_i Q_{i,j}^{+2})^{n_{i,j}} \right) \times \prod_{j=1}^J \left(\sum_{i=1}^I \sigma_i Q_{i,j}^{+2} \right)^{-n^j} \times \pi(\boldsymbol{\sigma}, \mathbf{Q}) \quad (7)$$

$$\propto \int \pi(\boldsymbol{\sigma}, \mathbf{Q}) \prod_{j=1}^J \left\{ \left(\prod_{i=1}^I (\sigma_i Q_{i,j}^{+2})^{n_{i,j}} \right) \frac{T_j^{n^j-1} \exp(-T_j \sum_i \sigma_i Q_{i,j}^{+2})}{\Gamma(n^j)} \right\} d\mathbf{T}, \quad (8)$$

where π is the prior. In order to obtain approximate $(\boldsymbol{\sigma}, \mathbf{Q})$ sampling we specify a Gibbs sampler for $(\boldsymbol{\sigma}, \mathbf{Q}, \mathbf{T})$ with target distribution

$$p(\boldsymbol{\sigma}, \mathbf{Q}, \mathbf{T}|\mathbf{n}) \propto \pi(\boldsymbol{\sigma}, \mathbf{Q}) \prod_{j=1}^J \left\{ \left(\prod_{i=1}^I (\sigma_i Q_{i,j}^{+2})^{n_{i,j}} \right) \frac{T_j^{n^j-1} \exp(-T_j \sum_i \sigma_i Q_{i,j}^{+2})}{\Gamma(n^j)} \right\}. \quad (9)$$

The sampler iterates the following steps:

[Step 1] Sample T_j independently, one for each biological sample $j = 1, \dots, J$,

$$T_j | \mathbf{Q}, \boldsymbol{\sigma}, \mathbf{n} \sim \text{Gamma}(n^j, \sum_i \sigma_i Q_{i,j}^{+2}).$$

[Step 2] Sample \mathbf{Q}_i independently, one for each OTU $i = 1, \dots, I$. The conditional density of $\mathbf{Q}_i = (Q_{i,1} \dots Q_{i,J})$ given $\boldsymbol{\sigma}, \mathbf{T}, \mathbf{n}$ is log-concave, and the random vectors \mathbf{Q}_i , $i = 1, \dots, I$, given $\boldsymbol{\sigma}, \mathbf{T}, \mathbf{n}$ are conditionally independent.

We simulate, for $j = 1, \dots, J$, from

$$p(Q_{i,j} | \mathbf{Q}_{i,-j}, \boldsymbol{\sigma}, \mathbf{T}, \mathbf{n}) \propto Q_{i,j}^{+2n_{i,j}} \times \exp(-T_j \sigma_i Q_{i,j}^{+2}) \times \exp\left(-\frac{(Q_{i,j} - \mu_{i,j})^2}{2s_j^2}\right), \quad (10)$$

where $\mathbf{Q}_{i,-j} = (Q_{i,1}, \dots, Q_{i,j-1}, Q_{i,j+1}, \dots, Q_{i,J})$, $\mu_{i,j} = E[Q_{i,j} | \mathbf{Q}_{i,-j}]$, $s_j^2 = \text{var}[Q_{i,j} | \mathbf{Q}_{i,-j}]$, with the proviso $0^0 = 1$. Since \mathbf{Q}_i is a multivariate Normal both $\mu_{i,j}$ and s_j have simple closed form expressions.

When $n_{i,j} = 0$ the density in (10) reduces to a mixture of truncated normals:

$$(1 - p_1) N(Q_{i,j}; \frac{\mu_{i,j}}{\Delta_{i,j}}, \frac{s_j^2}{\Delta_{i,j}}) I(Q_{i,j} > 0) + p_1 N(Q_{i,j}; \mu_{i,j}, s_j^2) I(Q_{i,j} \leq 0),$$

$$p_1 = \frac{\Phi(0; \mu_{i,j}, s_j^2) N(0; \frac{\mu_{i,j}}{\Delta_{i,j}}, \frac{s_j^2}{\Delta_{i,j}})}{\Phi(0; \mu_{i,j}, s_j^2) N(0; \frac{\mu_{i,j}}{\Delta_{i,j}}, \frac{s_j^2}{\Delta_{i,j}}) + N(0; \mu_{i,j}, s_j^2) \left(1 - \Phi(0; \frac{\mu_{i,j}}{\Delta_{i,j}}, \frac{s_j^2}{\Delta_{i,j}})\right)},$$

and $\Delta_{i,j} = 1 + 2\sigma_i T_j s_j^2$. Here $N(\cdot; \mu, s^2)$ and $\Phi(\cdot; \mu, s^2)$ are the density and cumulative density functions of a Normal variable with mean μ and variance s^2 .

When $n_{i,j} > 0$ the density $p[Q_{i,j} | \mathbf{Q}_{i,-j}, \boldsymbol{\sigma}, \mathbf{T}, \mathbf{n}]$ remains log-concave and the support becomes $(0, +\infty)$. We update $Q_{i,j}$ using a Metropolis-Hastings step with proposal

identical to the Laplace approximation $N(\hat{\mu}_{i,j}, \hat{s}_{i,j}^2)$ of the density in (10),

$$\hat{\mu}_{i,j} = \frac{\mu_{i,j}/s_j^2 + \sqrt{\mu_{i,j}^2/s_j^4 + 8n_{i,j}(2\sigma_i T_j + 1/s_j^2)}}{2(2\sigma_i T_j + 1/s_j^2)}, \quad \hat{s}_{i,j}^2 = \left(\frac{2n_{i,j}}{\hat{\mu}_{i,j}^2} + 2T_j\sigma_i + \frac{1}{s_j^2} \right)^{-1}. \quad (11)$$

Here $\hat{\mu}_{i,j}$ maximizes the density (10) and $\hat{s}_{i,j}^2$ is obtained from the second derivative of the log-density at $\hat{\mu}_{i,j}$. We found the approximation accurate. In Appendix D we provide bounds of the total variation distance between the target (10) and the approximation (11). When $n_{i,j}$ increases, the bound of the total variation decreases to zero. See also Figure A.1 in the Appendix.

[Step 3] Sample σ_i independently, one for each OTU $i = 1, \dots, I$, from the density $p(\sigma_i | \mathbf{Q}, \mathbf{T}, \mathbf{n}) \propto \pi(\sigma_i) \sigma_i^{n_i} \exp(-\sigma_i \sum_{j=1}^J T_j Q_{i,j}^{+2})$. The σ_i 's are a priori independent $\text{Beta}(\alpha/I, 1/2 - \alpha/I)$ variables. We use piecewise constant bounds for $\sigma \rightarrow \exp(-\sigma_i \sum_{j=1}^J T_j Q_{i,j}^{+2})$, $\sigma \in [0, 1]$, and an accept/reject step to sample from $p(\sigma_i | \mathbf{Q}, \mathbf{T}, \mathbf{n})$.

We now consider inference on Σ using the prior on \mathbf{Y} in subsection 2.3. The goal is to generate approximate samples of \mathbf{Y} from the posterior. We exploit the identity of the conditional distributions of \mathbf{Y} given $(\sigma, \mathbf{T}, \mathbf{Q}, \mathbf{n})$ and \mathbf{Q} . In order to sample \mathbf{Y} from the posterior we can therefore directly apply the MCMC transitions in Bhattacharya and Dunson (2011), with \mathbf{Q} replacing the observable variables in their work.

3.1 Self-consistent estimates of biological samples' similarity

We discuss an EM-type algorithm to estimate the correlation matrix \mathbf{S} of the vectors $(Q_{i,1}, \dots, Q_{i,J})$, $i = 1, \dots, I$. Under our construction in subsection 2.3) we interpret \mathbf{S} as the normalized Gram matrix ϕ between biological samples. In this subsection we describe an alternative procedure, distinct from the Gibbs sampler, which does not require tuning of the prior probability model. The algorithm can be used for MCMC initialization and for exploratory data analyses. It assumes that the observed OTU abundances are representative of the microbial distributions, i.e. $P^j = (n_{1,j}/n^j, \dots, n_{I,j}/n^j)$. Under this assumption, for each biological sample j ,

$$\begin{aligned} \sigma_i Q_{i,j}^{+2} \times I(n_{i,j} > 0) &\propto n_{i,j}, \quad i = 1, \dots, I, \\ \text{and } Q_{i,j} &\leq 0 \quad \text{when } n_{i,j} = 0. \end{aligned} \quad (12)$$

For σ_i , $i = 1, \dots, I$, we use a moment estimate $\hat{\sigma}_i = (1/J) \sum_j \left(n_{i',j} / \sum_{i' \neq i} n_{i',j} \right)$. The procedure uses these estimates and at iteration $t + 1$ generates the following results: **[Expectation]** Impute repeatedly \mathbf{Q} , $\ell = 1, \dots, D$ times, consistently with the constraints (12) and using a $N(0, \Sigma_t)$ joint distribution. Here Σ_t is the estimate of Σ , the covariance matrix of $(Q_{i,1}, \dots, Q_{i,J})$, after the t -th iteration. For each replicate $\ell = 1, \dots, D$, we fix $Q_{i,j}^\ell$ for all (i, j) pairs with strictly positive $n_{i,j}$ counts at $\sqrt{n_{i,j}/\hat{\sigma}_i}$, and sample jointly, conditional on these values, negative $Q_{i,j}^\ell$ values for the remaining (i, j) pairs with $n_{i,j} = 0$. We use these $Q_{i,j}^\ell$ values to approximate $\mathcal{L}(\Sigma)$, the full data log-likelihood, our target function as in any other EM algorithm.

[**Maximization**] Set Σ_{t+1} equal to the empirical covariance matrix of the $(Q_{i,1}^\ell, \dots, Q_{i,J}^\ell)$ vectors, thus maximizing the $\mathcal{L}(\Sigma)$ approximation.

We iterate until convergence of Σ_t . Then, after the last iteration, the inferred covariance matrix of $(Q_{i,1}, \dots, Q_{i,J})$ directly identifies an estimate of \mathbf{S} . We evaluated the algorithm using in-silico datasets from the simulation study in Section 5. Overall it generates estimates that are slightly less accurate compared to posterior estimation based on MCMC simulations. We use the datasets considered in Figure 2(a), with number of factors fixed at three and n^j at 100,000, for a representative example. In this case the average RV-coefficient between the true \mathbf{S} and the estimated matrix is 0.93 for the EM-type algorithm, and 0.95 for posterior simulations. In our work the described procedure reduced the computing time to approximately 10% compared to the Gibbs sampler. More details on this procedure are provided in Appendix E.

4 Visualizing uncertainty in ordination plots

Ordination methods such as Multidimensional Scaling of ecological distances or Canonical Correspondence Analysis are central in microbiome research. Given posterior samples of the model parameters, we use a procedure to plot credible regions in visualizations such as Fig 2(f). The methods that we consider here are all related to PCA and use the normalized Gram matrix \mathbf{S} between biological samples. We recall that in our model \mathbf{S} is the correlation matrix of $(Q_{i,1}, \dots, Q_{i,J})$. Based on a single posterior instance of \mathbf{S} , we can visualize biological samples in a lower dimensional space through PCA, with each biological sample projected once. Naively, one could think that simply overlaying projections of the principal component loadings generated from different posterior samples of \mathbf{S} on the same graph would show the variability of the projections. However, these super-impositions could be spurious if we carry out PCA for each \mathbf{S} sample separately. One possible problem is PC switching, when two PCs have similar eigenvalues. Another problem is the ambiguity of signs in PCA, which would lead to random signs of the loadings that result in symmetric groups of projections of the same biological sample at different sides of the axes. More generally PCA projections from different posterior samples of \mathbf{S} are difficult to compare, as the different lower dimensional spaces are not aligned.

We alternatively identify a consensus lower dimensional space for all posterior samples of \mathbf{S} (Escoufier, 1973; Lavit et al., 1994; Abdi et al., 2005). We list the three main steps used to visualize the variability of \mathbf{S} .

1. Identify a Gram matrix \mathbf{S}_0 that best summarizes K posterior samples of Gram matrix $\mathbf{S}_1, \dots, \mathbf{S}_K$. One simple criterion is to minimize L_2 loss element-wise. This leads to $\mathbf{S}_0 = (\sum_i \mathbf{S}_i)/K$. Alternatively, we can define \mathbf{S}_0 as the Gram matrix that maximizes similarity with $\mathbf{S}_1, \dots, \mathbf{S}_K$. One possible similarity metric between two symmetric square matrices \mathbf{A} and \mathbf{B} is the RV-coefficient (Robert and Escoufier, 1976), $\text{RV}(\mathbf{A}, \mathbf{B}) = \text{Tr}(\mathbf{AB})/\sqrt{\text{Tr}(\mathbf{AA})\text{Tr}(\mathbf{BB})}$. We refer to Holmes (2008) for a discussion on RV-coefficients.

2. Identify the lower dimensional consensus space V based on \mathbf{S}_0 . Assume we want $\dim(V) = 2$; the basis of V will be the orthonormal eigenvectors \mathbf{v}_1 and \mathbf{v}_2 of \mathbf{S}_0 corresponding to the largest eigenvalues λ_1 and λ_2 . The configuration of all biological samples in V is visualized by projecting rows of \mathbf{S}_0 onto V : $(\boldsymbol{\psi}_1^0, \boldsymbol{\psi}_2^0) = \mathbf{S}_0(\mathbf{v}_1\lambda_1^{-1/2}, \mathbf{v}_2\lambda_2^{-1/2})$. As in a standard PCA, this configuration best approximates the Gram matrix in the L_2 sense: $(\boldsymbol{\psi}_1^0, \boldsymbol{\psi}_2^0) = \operatorname{argmin}_{\langle \boldsymbol{\psi}_1, \boldsymbol{\psi}_2 \rangle = 0} \|\mathbf{S}_0 - (\boldsymbol{\psi}_1, \boldsymbol{\psi}_2)(\boldsymbol{\psi}_1, \boldsymbol{\psi}_2)'\|^2$.
3. Project the rows of posterior sample \mathbf{S}_k onto V by $(\boldsymbol{\psi}_1^k, \boldsymbol{\psi}_2^k) = \mathbf{S}_k(\mathbf{v}_1\lambda_1^{-1/2}, \mathbf{v}_2\lambda_2^{-1/2})$. Overlaying all the $\boldsymbol{\psi}^k$ displays uncertainty of \mathbf{S} in the same linear subspace. Posterior variability of the biological samples' projections is visualized in V by plotting each row of the matrices $(\boldsymbol{\psi}_1^k, \boldsymbol{\psi}_2^k)$, $k = 1, \dots, K$, in the same figure. A contour plot is produced for each biological sample (see for example Fig 2(f)) to facilitate visualization of the posterior variability of its position in the consensus space V .

5 Simulation Study

In this section, we evaluate the procedure described in Section 3 and explore whether the shrinkage prior allows us to infer the number of factors and the normalized Gram matrix between biological samples \mathbf{S} . We also consider the estimates $E(P^j|\mathbf{n})$ obtained with our joint model, one for each biological sample j , and compare their precision with the empirical estimator.

We first defined a scenario with distributions P^j generated from the prior (1), with $I = 68$ OTUs and $J = 22$ biological samples. The true number of factors is L , and for biological samples $j = 1, \dots, J/2$, the vector $\mathbf{Y}^j = (Y_{l,j}, 1 \leq l \leq L)$ has elements $l = L/2 + 1, \dots, L$ equal to zero, while symmetrically, for $j = J/2 + 1, \dots, J$, the vectors \mathbf{Y}^j have the elements $l = 1, \dots, L/2$ equal to zero. The underlying normalized Gram matrix \mathbf{S} is therefore block-diagonal. After generating the distributions P^j , we sampled with fixed total counts (n^j) per biological sample $n^j = 1,000$. We produced 50 replicates with $L = 3, 6$ and 9 . In our simulations the non-zero components $Y_{l,j}$'s are independent standard Normal.

We use PCA-type summaries for the posterior samples of \mathbf{Y} generated from $p(\mathbf{Y}|\mathbf{n})$. Computations are based on the $J \times J$ normalized Gram matrix \mathbf{S} . At each MCMC iteration we generate approximate samples \mathbf{Y} from the posterior, compute \mathbf{S} by normalizing the Gram matrix $\mathbf{Y}'\mathbf{Y}$, and operate standard spectral decomposition on \mathbf{S} . This allows us to estimate the ranked eigenvalues, i.e. the principal components' variance of our \mathbf{Q} latent vectors (after normalization), by averaging over the MCMC iterations. Figure 2(a) shows the variability captured by the first 10 principal components, with the box-plots illustrating posterior means' variability across our 50 replicates. The proportion of variability associated to each *principal component* decreases rapidly after the *true* number of factors $L = 3, 6, 9$. This suggests that the shrinkage model (Bhattacharya and Dunson, 2011) tends to produce posterior distributions for our \mathbf{Y} latent variables that concentrates around a linear subspace.

Figure 2(c) illustrates the accuracy of the estimated normalized Gram matrix $\hat{\mathbf{S}}$ with n^j equal to 1,000, 10,000 and 100,000. We estimated the unknown $J \times J$ normalized Gram matrix \mathbf{S} with the posterior mean of the normalized Gram matrix, which we approximate by averaging over MCMC iterations. We summarized the accuracy using the RV coefficient between $\hat{\mathbf{S}}$ and \mathbf{S} , see Robert and Escoufier (1976) for a discussion on this metric. The box-plots illustrate variability of estimates' accuracy across 50 simulation replicates. As expected, when the total counts per sample increases from 10,000 to 100,000, we only observe limited gain in accuracy. Indeed the overall number of observed OTUs with positive counts per biological sample remains comparable, with expected values equal to 30 and 33 when the total counts per biological sample are fixed at 10,000 and 100,000 respectively. We also note that when L increases, the accuracy decreases.

We investigate interpretability of our model by using distributions P^j generated from a probability model that slightly differs from the prior. More precisely, the i th random weight in P^j , conditionally on \mathbf{Y} and \mathbf{X} , is defined proportional to a monotone function of $\langle \mathbf{X}_i, \mathbf{Y}^j \rangle^+$. We considered for example

$$P^j(A) = \frac{\sum_i \sigma_i \langle \mathbf{X}_i, \mathbf{Y}^j \rangle^{+a} I_{Z_i}(A)}{\sum_i \sigma_i \langle \mathbf{X}_i, \mathbf{Y}^j \rangle^+}, \quad a > 0. \quad (13)$$

When the monotone function is quadratic the probability model becomes identical to our prior. In Figure 2(b) and Figure 2(d) we used model (13) with $a = 1$ to generate datasets. We repeated the same simulation study summarized in the previous paragraphs.

We evaluated the effectiveness of borrowing information across biological samples for estimating the vectors P^j . The accuracy metric that we used is the total variation distance. We compared the Bayesian estimator $E(P^j | \mathbf{n})$ and the empirical estimator \tilde{P}^j which assigns mass $n_{i,j}/n^j$ to the i^{th} OTUs. The advantage of pooling information varies with the similarity between biological samples. To reflect this, we generated P^j with non-zero components of \mathbf{Y} sampled from a zero mean multivariate Normal with $\text{cov}(Y_{l,j}, Y_{l,j'})$ equal to θ . We considered the case when P^j is generated either from our prior or model (13) with $a = 0.5, 1, 3$. In addition, we considered $\theta = 0.5, 0.75, 0.95$, $I = 68, J = 22$ and $L = 3$, while n^j varies from 10 to 100.

The results are summarized in Figure 2(e) which shows the average difference in total variation, contrasting the Bayesian and empirical estimators. The results both, when the model is correctly specified and mis-specified, quantify the advantages in using a joint Bayesian model.

We complete this section with one illustration of the method in Section 4. We simulate a dataset with two clusters by generating $Y_{l,j}$ for $l = 1, \dots, L$ from $N(-3, 1)$ when $j = 1, \dots, J/2$ and from $N(3, 1)$ when $j = J/2 + 1, \dots, J$. All $Y_{l,j}$ are different from zero. We expected a low n^j to be sufficient for detecting the clusters. We sampled P^j from the prior and set $J = 22, I = 68, L = 3$ and $n^j = 100$. The PC plot and the biological sample specific confidence regions are shown in Figure 2(f). In the PC plot the two clusters are illustrated with different colors. In this simulation exercise the posterior credible regions leave little ambiguity both on the presence of clusters and

also on samples-specific cluster membership.

6 Application to microbiome datasets

In this section, we apply our Bayesian analysis to two microbiome datasets. We show that our method gives results that are consistent with previous studies and we show our novel visualization of uncertainty in ordination plots. We start with the Global Patterns data (Caporaso et al., 2011) where human-derived and environmental biological samples are included. We then considered data on the vaginal microbiome (Ravel et al., 2011).

6.1 GlobalPatterns dataset

The GlobalPatterns dataset includes 26 biological samples derived from both human and environmental specimens. There are a total of 19,216 OTUs and the average total counts per biological sample is larger than 100,000. We collapsed all taxa OTUs to the genus level –a standard operation in microbiome studies– and yielded 996 distinct genera. We treated these genera as OTUs’ and fit our model to this collapsed dataset. We ran one MCMC chain for 50,000 iterations and recorded posterior samples every 10 iterations.

We first performed a clustering analysis of biological samples using Partitioning Among Medoids (PAM) (Hastie et al., 2003). For each posterior sample of the model parameters, we computed $P^j|\mathbf{n}$ for $j = 1, \dots, J$ and calculated the Bray-Curtis dissimilarity matrix between biological samples. By averaging over the MCMC iterations and the clustering results from each dissimilarity matrix, we obtained the posterior probability of two biological samples being clustered together. Figure 3(a) illustrates the clustering probabilities. We can see that biological samples belonging to a specific specimen type are tightly clustered together while different specimens tend to define separate clusters. This is consistent with the conclusion in Caporaso et al. (2011), where the authors suggest, that within specimen microbiome variations are limited when compared to variations across specimen types. We also observed that biological samples from the skin are clustered with those from the tongue. This is to some extent an expected result, because both specimens are derived from humans, and because the skin microbiome has often OTUs frequencies comparable to other body sites (Grice and Segre, 2011).

We then visualized the biological samples using ordination plots and applying the method described in Section 4. We fixed the dimension of the consensus space V at three. We plotted all biological samples’ projections onto V along with contours to visualize their posterior variability. The results are shown in Figure 3(b-d). We observe a clear separation between human-derived (tongue, skin and feces) biological samples and biological samples from free environments. This separation is mostly identified by the first two compromise axes. The third axis defines a saline/non-saline samples separation. Biological samples derived from saline environment (e.g. Ocean) are well separated when projected on this axis from those derived from non-saline

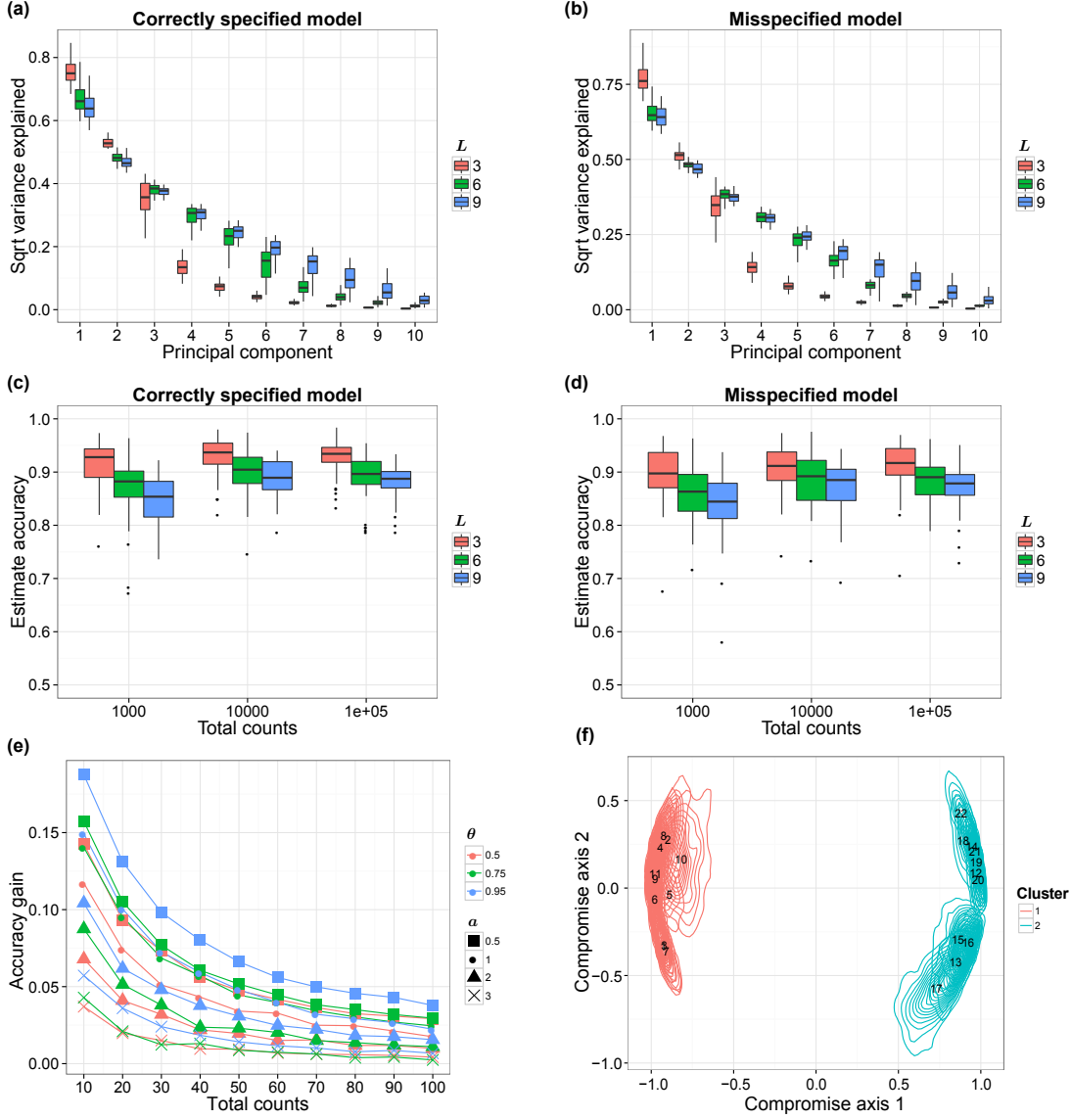


Figure 2: (a-b) Estimated proportion of variability captured by the first 10 PCs. Each box-plot here shows the variability of the estimated proportion across 50 simulation replicates. We show the results when the data are generated from the prior (Panel a) and from the model in (13) with $a = 1$ (Panel b). (c-d) Accuracy of the correlation matrix estimates $\hat{\mathbf{S}}$. The box-plots show the variability of the accuracy in 50 simulation replicates, with data generated from the prior (Panel c) and from model (13) with $a = 1$ (Panel d). We vary the number of factors L (colors) and n^j and show the corresponding accuracy variations. (e) Comparison between Bayesian estimates of the underlying microbial distributions P^j and the empirical estimates. We consider the average total variation difference, averaging across all J biological samples. Each curve shows the relationship between n^j and average accuracy gain. We set $L = 3$ and the parameter a varies from 0.5 to 3 (shapes). The similarity parameter θ is equal to 0.5, 0.75 or 0.95 (colors). (f) PCoA plot with confidence regions. We visualize the confidence regions using the method in Section 4. Each contour illustrates the uncertainty of a single biological sample's position. Colors indicate cluster membership and annotated numbers are biological samples' IDs.

environment (e.g. Creek freshwater). We observed small 95% credible regions for all biological samples projections. This low level of uncertainty captured by the small credible regions in Figure 3(b-d) is mainly explained by the large total counts n^j for all biological samples.

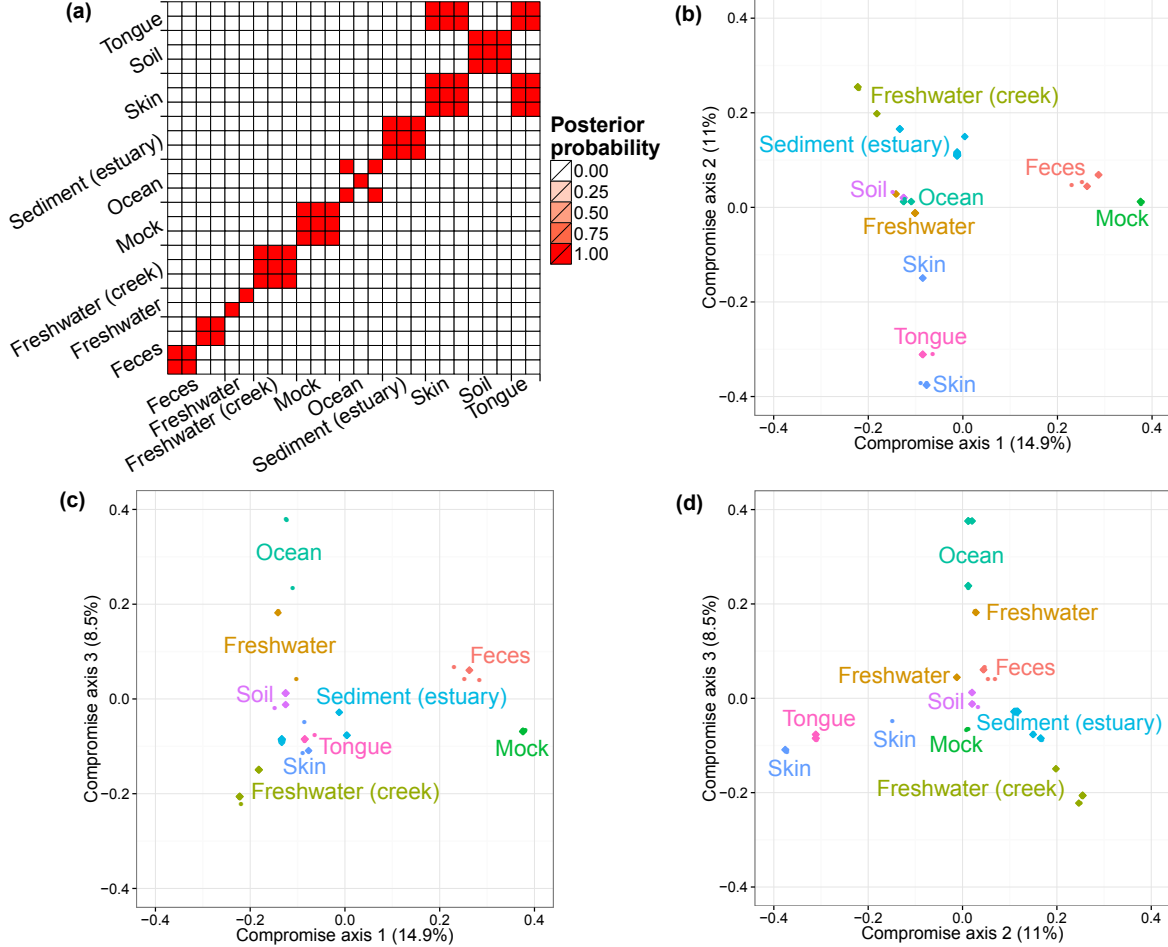


Figure 3: (a) Posterior Probability of each pair of biological samples (j, j') being clustered together. The labels on axes indicate the environment of origin for each biological sample. (b-d) Ordination plots of biological samples and 95% posterior credible regions. We illustrate the first three compromise axes with three panels. The percentages on the three axes are the ratios of the corresponding \mathbf{S}_0 eigenvalues and the trace of the matrix. The credible regions for some biological samples are so small that appears as single points. Colors and annotated text indicate the environments.

6.2 The Vaginal Microbiome

We also consider a dataset previously presented in Ravel et al. (2011) which contains a larger number of biological samples (900) and a simpler bacterial community structure.

These biological samples are derived from 54 healthy women. Multiple biological samples are taken from each individual, ranging from one to 32 biological samples per individual. Each woman has been classified, before our microbiome sequencing data were generated, into vaginal community state subtypes (CST). This dataset contains only species level taxonomic information and we filtered OTUs by occurrence. We only retain species with more than 5 reads in at least 10% of biological samples. This filtering resulted in 31 distinct OTUs. We ran one MCMC chain with 50,000 iterations.

We performed the same analyses as in the previous subsection. The results are shown in Figure 4. Clustering probabilities indicate strong within CST similarity (panel a). There is one exception, CST IV-A samples, in some cases, presents low levels of similarities when compared to each other, and tend to cluster with CST I, CST III and CST IV-B samples. This is because CST IV-A is characterized as a highly heterogeneous subtype (Ravel et al., 2011). The ordination plots are consistent with the discoveries in Ravel et al. (2011). A tetrahedron shape is recovered and CST I, II, III, IV-B occupy the four vertices. CST II is well separated from other CSTs by the third axis. We also observed a sub-clustering in CST II which has not been detected and discussed in Ravel et al. (2011). This difference in the results can be due to distinct clustering metrics in the analyses.

Note that there are two biological samples with large credible regions, indicating high uncertainty of the corresponding positions. This uncertainty propagates on their cluster membership. Both biological samples have small total counts compared to the others. The lack of precision when using biological samples with small sequencing depth leads to high uncertainty in ordination and classification. It is therefore important to account for uncertainty in the validation of subgroups biological differences—in our case CST subtypes—based on microbiome profiling. Our example suggests also the importance of uncertainty summaries when microbiome profiles are used to classify samples. Uncertainty summaries allow us to retain all samples, including those with low counts, without the risk of overinterpreting the estimated locations and projections. This also argues for the retention of raw counts in microbiome studies (McMurdie and Holmes, 2014). By using raw counts, we can evaluate the uncertainty of our estimates and exploit the information and statistical power carried by the full dataset; whereas if we downsample the data we lose information and increase uncertainty on the projections.

It is ubiquitous to have biological samples with relevant differences in their total counts, and in some cases the number of OTUs and the total number of reads can be comparable. In this cases, the empirical estimates of microbial distributions are not reliable, and an assessment of the uncertainty is necessary for downstream analyses. The two biological samples with low total counts in the vaginal microbiome dataset are examples. For biological samples with a scarce amount of data our model provides measures of uncertainty and allows uncertainty visualizations with ordination plots.

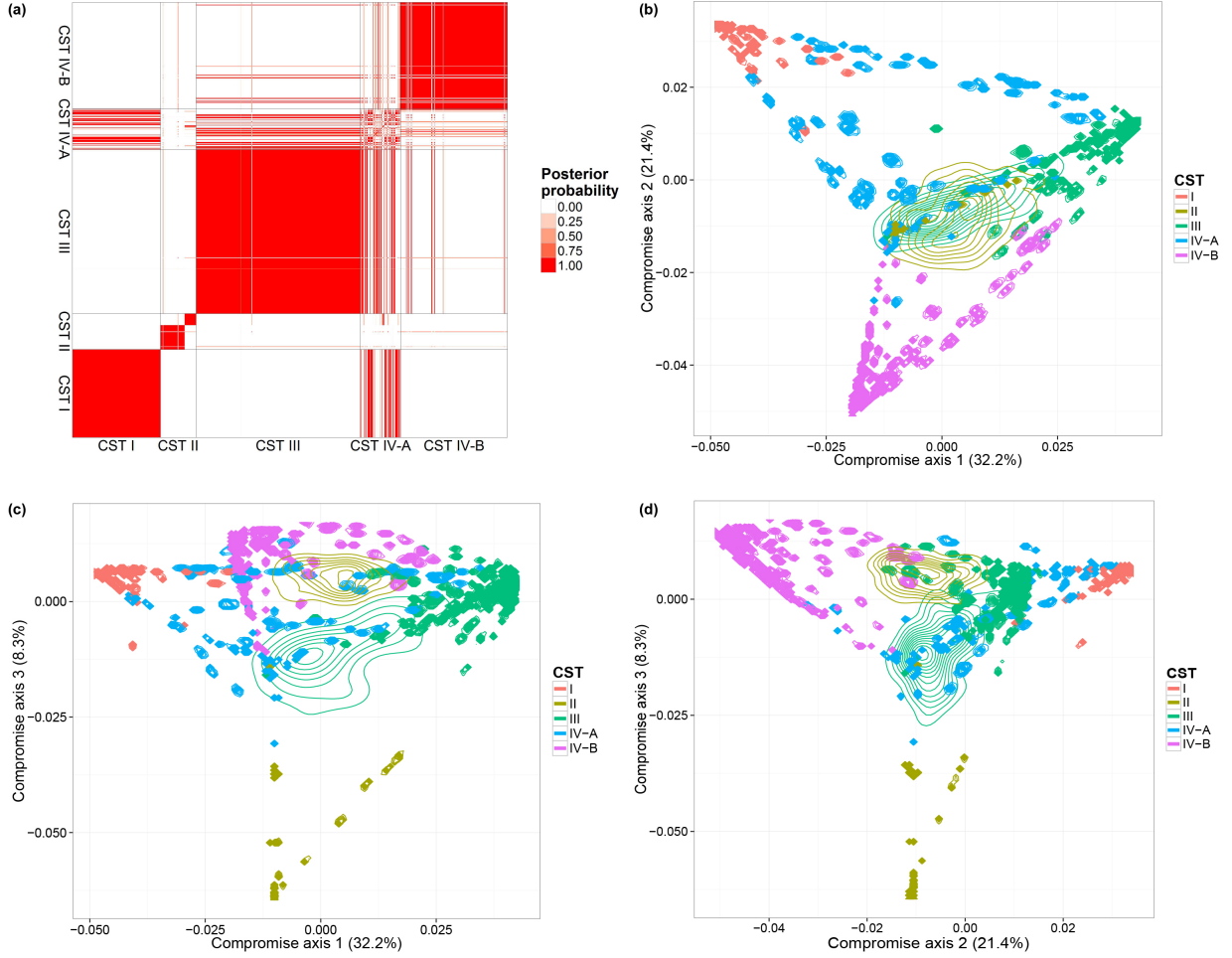


Figure 4: (a) Posterior Probability of each pair of biological samples (j, j') being clustered together. The labels on axes indicate the CST for each biological sample. (b-d) Ordination plots of biological samples and posterior credible regions. We illustrate the first three compromise axes with three panels. The percentages on the three axes are the ratios of the corresponding \mathbf{S}_0 eigenvalues and the trace of the matrix. Colors and indicate CSTs.

7 Conclusion

We propose a joint model for multinomial sampling of OTUs in multiple biological samples. We apply a prior from Bayesian factor analysis to estimate the similarity between biological samples, which is summarized by a Gram matrix. Simulation studies give evidence that this parameter is recovered by the Bayes estimate, and in particular, the inherent dimensionality of the latent factors is effectively learned from the data. The simulation also demonstrated that the analysis yields more accurate estimates of the microbial distributions by borrowing information across biological samples.

In addition, we provide a robust method to visualize the uncertainty in ecological ordinations, furnishing each point in the plot with a credible region. Two published microbiome datasets were analyzed and the results are consistent with previous findings. The second analysis demonstrates that the level of uncertainty can vary across biological samples due to differences in sampling depth, which underlines the importance of modeling multinomial sampling variations coherently. We believe our analysis will mitigate artifacts arising from rarefaction, thresholding of rare species, and other preprocessing steps.

There are several directions for development which are not explored here. We highlight the possibility of incorporating prior knowledge about the biological samples, such as the subject or group identifier in a clinical study. To model fixed effects, the latent factors \mathbf{Y}^j can be augmented by a vector of covariates $(b_1 w_1^j, \dots, b_p w_p^j)$, whose coefficients b could be given a normal prior, for example. The posterior distribution of the coefficients could be used to infer the magnitude of covariates' effects. A less straightforward extension involves moving away from the assumption of *a priori* exchangeability between OTUs to include prior information about phylogenetic or functional relationships between them. In our present analysis, these relationships are not taken into account in the definition of the prior for microbial distributions.

8 Acknowledgements

B. Ren is supported by National Science Foundation under Grant No. DMS-1042785. Stefano Favaro is supported by the European Research Council (ERC) through StG N-BNP 306406. L. Trippa has been supported by the Claudia Adams Barr Program in Innovative Basic Cancer Research. S. Holmes was supported by the NIH grant R01AI112401. We thank Persi Diaconis, Kris Sankaran and Lan Huong Nguyen for helpful suggestions and improvements.

A Approximating a Poisson Process using Beta random variable

Consider the approximation problem where we want to approximate the Poisson process on $(0, 1)$ with intensity measure given by the density with respect to Lebesgue

measure: $\nu(\sigma) = \alpha\sigma^{-1}(1-\sigma)^{-1/2}$. The finite counting process we derived is n iid samples drawn from Beta distribution $\text{Beta}(\epsilon_n, 1/2 - \epsilon_n)$ where $\epsilon_n < 1/2$ is a very small positive quantity and change with n . Denote the Poisson process as $N(t)$ and the approximated process as $N'_n(t)$, we can first calculate the probability of having m counts in interval $(\delta, t]$, where $m \leq n$, $t < 1$ and $0 < \delta \ll 1$:

$$P[N((\delta, t]) = m] = \frac{\left[\int_{\delta}^t \alpha\sigma^{-1}(1-\sigma)^{-1/2} d\sigma \right]^m}{m!} \exp\left(- \int_{\delta}^t \alpha\sigma^{-1}(1-\sigma)^{-1/2} d\sigma\right)$$

$$P[N'_n((\delta, t]) = m] = \binom{n}{m} \left(\frac{1}{\text{Beta}(\epsilon_n, 1/2 - \epsilon_n)} \int_{\delta}^t \sigma^{-1+\epsilon_n}(1-\sigma)^{-1/2-\epsilon_n} d\sigma \right)^m$$

$$\left(1 - \frac{1}{\text{Beta}(\epsilon_n, 1/2 - \epsilon_n)} \int_{\delta}^t \sigma^{-1+\epsilon_n}(1-\sigma)^{-1/2-\epsilon_n} d\sigma \right)^{n-m}$$

The corresponding moment generating functions are

$$M_N(\lambda) = \exp\left[(e^{\lambda} - 1) \int_{\delta}^t \alpha\sigma^{-1}(1-\sigma)^{-1/2} d\sigma \right]$$

$$M_{N'_n}(\lambda) = \left[\frac{e^{\lambda} - 1}{\text{Beta}(\epsilon_n, 1/2 - \epsilon_n)} \int_{\delta}^t \sigma^{-1+\epsilon_n}(1-\sigma)^{-1/2-\epsilon_n} d\sigma + 1 \right]^n$$

These two MGFs will be the same asymptotically when

$$\lim_{n \rightarrow \infty} \frac{n}{\text{Beta}(\epsilon_n, 1/2 - \epsilon_n)} \int_{\delta}^t \sigma^{-1+\epsilon_n}(1-\sigma)^{-1/2-\epsilon_n} d\sigma = \alpha \int_{\delta}^t \sigma^{-1}(1-\sigma)^{-1/2} d\sigma$$

This can be satisfied when $\epsilon_n = \alpha/n$. In this case,

$$\lim_{n \rightarrow \infty} \frac{n(\sigma/(1-\sigma))^{\epsilon_n}}{\text{Beta}(\epsilon_n, 1/2 - \epsilon_n)} = \alpha.$$

Also, when n is large enough, the map $n \mapsto \frac{n(\sigma/(1-\sigma))^{\epsilon_n}}{\text{Beta}(\epsilon_n, 1/2 - \epsilon_n)}$ is a non-increasing function. By Lebesgue's monotone convergence theorem, we can establish the convergence of the first integral to the second.

Using this result, we can construct the weak convergence of the finite dimension distribution: $(N'(\delta, t_1], \dots, N'(\delta, t_n]) \xrightarrow{d} (N(\delta, t_1], \dots, N(\delta, t_n])$. This follows by a direct application of the multinomial theorem.

Now we need to verify the tightness condition, this is automatically satisfied as the process defined in terms of Beta random variables is a càdlàg process (Daley and Vere-Jones, 1988) (Theorem 11.1. VII and Proposition 11.1. VIII, iv, Volume 2). Therefore we prove the weak convergence of the process built by independent Beta random variables.

B Proof of Proposition 1

We use the notation $P^j(\cdot) = \frac{\sum_i I(Z_i \in \cdot) \sigma_i Q_{i,j}^{+2}}{\sum_i \sigma_i Q_{i,j}^{+2}}$ where $Q_{i,j} = \langle \mathbf{X}_i, \mathbf{Y}^j \rangle$. Denote $((Q_{i,j}, Q_{i,j'}), i \geq 1)$ as \mathbf{Q} . The joint distribution of $(Q_{i,j}, Q_{i,j'})$ is a multivariate normal with mean $\mathbf{0}$ and

covariance $\phi(j, j')$, and the vectors $(Q_{k,j}, Q_{k,j'})$, $k = 1, 2, \dots$, are independent. We derive an expression for the covariance

$$\begin{aligned} \text{cov}[P^j(A), P^{j'}(A)] &= E[E[P^j(A)P^{j'}(A)|\sigma, \mathbf{Q}]] - E[P^j(A)]E[P^{j'}(A)] \\ &= (G(A) - G^2(A))E \left[\frac{\sum_i \sigma_i^2 Q_{i,j}^{+2} Q_{i,j'}^{+2}}{\sum_i \sigma_i Q_{i,j}^{+2} \sum_k \sigma_k Q_{k,j'}^{+2}} \right]. \end{aligned}$$

The variance is

$$\text{var}[P^j(A)] = (G(A) - G^2(A))E \left[\frac{\sum_i \sigma_i^2 Q_{i,j}^{+4}}{\sum_i \sigma_i Q_{i,j}^{+2} \sum_k \sigma_k Q_{k,j}^{+2}} \right].$$

It follows that

$$\text{corr}[P^j(A), P^{j'}(A)] = E \left[\frac{\sum_i \sigma_i^2 Q_{i,j}^{+2} Q_{i,j'}^{+2}}{\sum_i \sigma_i Q_{i,j}^{+2} \sum_k \sigma_k Q_{k,j'}^{+2}} \right] \times \left(E \left[\frac{\sum_i \sigma_i^2 Q_{i,j}^{+4}}{\sum_i \sigma_i Q_{i,j}^{+2} \sum_k \sigma_k Q_{k,j}^{+2}} \right] \right)^{-1}.$$

The correlation is independent of the set A .

C Proof of Proposition 2

We follow the framework of proofs for Theorem 1 and Theorem 3 in Barrientos et al. (2012). Let $\mathcal{P}(\mathcal{Z})$ be the set of all Borel probability measures defined on $(\mathcal{Z}, \mathcal{F})$ and $\mathcal{P}(\mathcal{Z})^d$ the product space of d $\mathcal{P}(\mathcal{Z})$. Assume $\Theta \subset \mathcal{Z}$ is the support of G . To show the prior assigns strictly positive probability to the neighbourhood in Proposition 2, it is sufficient to show such neighbourhood contains certain subset-neighbourhoods with positive probability. As in Barrientos et al. (2012), we consider the subset-neighbourhoods U :

$$U(G_1, \dots, G_d, \{A_{i,j}\}, \epsilon^*) = \prod_{i=1}^d \{F_i \in \mathcal{P}(\Theta) : |F_i(A_{i,j}) - G_i(A_{i,j})| < \epsilon^*, j = 1, \dots, m_i\},$$

where G_i is a probability measure absolutely continuous w.r.t. G for $i = 1, \dots, d$, $A_{i,1}, \dots, A_{i,m_i} \subset \Theta$ are measurable sets with G_i -null boundary and $\epsilon^* > 0$. The existence of such subset-neighbourhoods is proved in Barrientos et al. (2012). We then define sets $B_{\nu_{1,1}, \nu_{m_d,d}}$ for each $\nu_{i,j} \in \{0, 1\}$ as

$$B_{\nu_{1,1}, \nu_{m_d,d}} = \bigcap_{i=1}^d \bigcap_{j=1}^{m_i} A_{i,j}^{\nu_{i,j}},$$

where $A_{i,j}^1 = A_{i,j}$ and $A_{i,j}^0 = A_{i,j}^c$. Set

$$J_\nu = \{\nu_{1,1}, \nu_{m_d,d} : G(B_{\nu_{1,1}, \nu_{m_d,d}}) > 0\},$$

and let \mathcal{M} be a bijective mapping from J_ν to $\{0, \dots, k\}$ where $k = |J_\nu| - 1$. Therefore we can simplify the notation using $A_{\mathcal{M}(\nu)} = B_\nu$ for every $\nu \in J_\nu$. Define a vector $\mathbf{s}_i = (w_{i,0}, \dots, w_{i,k}) = (Q_i(A_0), \dots, Q_i(A_k))$ that belongs to the k -simplex Δ_k . Set

$$B(\mathbf{s}_i, \epsilon) = \{(w_0, \dots, w_k) \in \Delta_k : |Q_i(A_j) - w_j| < \epsilon, j = 0, \dots, k\},$$

where $\epsilon = 2^{-\sum_{i=1}^d m_i} \epsilon^*$. The derivation in Barrientos et al. (2012) suggests a sufficient condition for assigning positive mass to $U(G_1, \dots, G_d, \{A_{i,j}\}, \epsilon^*)$ is

$$\Pi([P^i(A_0), \dots, P^i(A_k)] \in B(\mathbf{s}_i, \epsilon), i = 1, \dots, d) > 0. \quad (\text{A.1})$$

Here Π is the prior.

Now consider the following conditions

$$\text{C.1 } w_{i,l} - \epsilon_0 < \sigma_{l+1} Q_{l+1,i}^{+2} < w_{i,l} + \epsilon_0 \text{ for } i = 1, \dots, d \text{ and } l = 0, \dots, k.$$

$$\text{C.2 } 0 < \sum_{l>k+1} \sigma_l Q_{l,i}^{+2} < \epsilon_0.$$

$$\text{C.3 } Z_{l+1} \in A_l \text{ for } l = 0, \dots, k.$$

ϵ_0 in the above conditions satisfies the following inequality

$$\begin{aligned} \frac{w_{(i,l)} - \epsilon_0}{1 + (k+2)\epsilon_0} &\geq w_{(i,l)} - \epsilon \\ \frac{w_{(i,l)} + 2\epsilon_0}{1 - (k+1)\epsilon_0} &\leq w_{(i,l)} + \epsilon \end{aligned}$$

for $i = 1, \dots, d$ and $l = 0, \dots, k$. This system of inequalities can be satisfied when k is large enough. If conditions (C.1) to (C.3) hold, it follows that $[P^i(A_0), \dots, P^i(A_k)] \in B(\mathbf{s}_i, \epsilon)$ for $i = 1, \dots, d$. Therefore, we have

$$\begin{aligned} &\Pi([P^i(A_0), \dots, P^i(A_k)] \in B(\mathbf{s}_i, \epsilon), i = 1, \dots, d) \geq \\ &\prod_{l=0}^k \Pi(w_{(i,l)} - \epsilon_0 < \sigma_{l+1} Q_{l+1,i}^{+2} < w_{(i,l)} + \epsilon_0, i = 1, \dots, d) \times \\ &\Pi\left(\sum_{l>k+1} \sigma_l Q_{l,i}^{+2} < \epsilon_0, i = 1, \dots, d\right) \times \\ &\prod_{l=0}^k \Pi(Z_{l+1} \in A_l) \times \Pi(Z_l \in \mathcal{Z} \mid l = k+2, \dots). \end{aligned}$$

Since $(Q_{l,1}, \dots, Q_{l,d})$ are multivariate normal random vectors with strictly positive definite covariance matrix and σ_l are always positive, the vector $(\sigma_{l+1} Q_{l+1,i}^{+2}, i = 1, \dots, d)$ has full support on \mathbb{R}^{k+1} and will assign positive probability to any subset of the space. It follows that

$$\Pi(w_{i,l} - \epsilon_0 < \sigma_{l+1} Q_{l+1,i}^{+2} < w_{i,l} + \epsilon_0, i = 1, \dots, d) > 0 \text{ for } l = 0, \dots, k.$$

Using the Gamma process argument, we know $\sum_{l>k+1} \sigma_l Q_{l,i}^{+2}$ is a total mass for a well-defined Gamma process and thus will always be positive for all i . It follows that

$$\Pi\left(\sum_{l>k+1} \sigma_l Q_{l,i}^{+2} < \epsilon_0, i = 1, \dots, d\right) > 0.$$

Since \mathcal{Z} is the topological support of G , it follows that $P(Z_{i+1} \in A_i) > 0$ and $P(Z_i \in \mathcal{Z}) = 1$. Combining these facts, we prove that Equation (A.1) holds.

D Total variation bound of Laplace approximate in Equation (11)

We consider the class of densities $g(x; k, \mu, s^2)$

$$g(x; k, \mu, s^2) \propto I(x \geq 0) x^{2k} f(x; \mu, s^2), k \in \mathbb{N}^+$$

where $f(x; \mu, s^2)$ is the density function of $N(\mu, s^2)$. The Laplace approximation of $g(x; k, \mu, s^2)$ is written as $f(x; \hat{\mu}, \hat{s}^2)$. Here $\hat{\mu} = \operatorname{argmax}_x g(x; k, \mu, s^2)$ and $\hat{s}^2 = -((\partial^2 \log(g)/\partial x^2)|_{\hat{\mu}})^{-1}$. We want to calculate the total variation distance between density $f(x; \hat{\mu}, \hat{s}^2)$ and $g(x; k, \mu, s^2)$, denoted as $d_{TV}(f(x; \hat{\mu}, \hat{s}^2), g(x; k, \mu, s^2))$.

Define class of functions $V(x; k, \mu)$ for $k \in \mathbb{N}^+, \mu > 0$:

$$V(x; k, \mu) = \begin{cases} 2k [\log(x/\mu) - (x/\mu - 1) + \frac{1}{2}(x/\mu - 1)^2] & x > 0 \\ -\infty & x \leq 0 \end{cases}$$

This function is non-decreasing and when $x = \mu$, $V(x; k, \mu) = 0$, $dV/dx = 0$ and $d^2V/dx^2 = 0$.

It follows that

$$\log g(x; k, \mu, s^2) - \log f(x; \hat{\mu}, \hat{s}^2) = V(x; k, \hat{\mu}) + a_0 + a_1 x + a_2 x^2.$$

Moreover, since the $\hat{\mu}$ is the mode of both $g(x; k, \mu, s^2)$ and $f(x; \hat{\mu}, \hat{s}^2)$, and the second derivative of $\log g(x; k, \mu, s^2)$ and $\log f(x; \hat{\mu}, \hat{s}^2)$ are identical at $x = \hat{\mu}$, we can find that $a_1 = a_2 = 0$. Hence,

$$\log g(x; k, \mu, s^2) - \log f(x; \hat{\mu}, \hat{s}^2) = V(x; k, \hat{\mu}) + a_0$$

and $g(x; k, \mu, s^2) = \exp(V(x; k, \hat{\mu}) + a_0) f(x; \hat{\mu}, \hat{s}^2)$.

Since $V(x; k, \hat{\mu})$ is monotone increasing, the total variation distance between $g(x; k, \mu, s^2)$ and $f(x; \hat{\mu}, \hat{s}^2)$ can be expressed as

$$\begin{aligned} d_{TV}(g(x; k, \mu, s^2), f(x; \hat{\mu}, \hat{s}^2)) &= \int_{x_0}^{+\infty} [\exp(V(x; k, \hat{\mu}) + a_0) - 1] f(x; \hat{\mu}, \hat{s}^2) dx \\ &= \int_{-\infty}^{x_0} [1 - \exp(V(x; k, \hat{\mu}) + a_0)] f(x; \hat{\mu}, \hat{s}^2) dx \end{aligned}$$

where $V(x_0; k, \hat{\mu}) = -a_0$. If $a_0 \leq 0$, we have $x_0 \geq \hat{\mu}$ and

$$\begin{aligned} & \int_{x_0}^{+\infty} [\exp(V(x; k, \hat{\mu}) + a_0) - 1] f(x; \hat{\mu}, \hat{s}^2) dx \\ & \leq \int_{x_0}^{+\infty} [\exp(V(x; k, \hat{\mu})) - 1] f(x; \hat{\mu}, \hat{s}^2) dx \\ & \leq \int_{\hat{\mu}}^{+\infty} [\exp(V(x; k, \hat{\mu})) - 1] f(x; \hat{\mu}, \hat{s}^2) dx \end{aligned}$$

Similarly, if $a_0 \geq 0$, we have

$$\int_{-\infty}^{x_0} [1 - \exp(V(x; k, \hat{\mu}) + a_0)] f(x; \hat{\mu}, \hat{s}^2) dx \leq \int_{-\infty}^{\hat{\mu}} [1 - \exp(V(x; k, \hat{\mu}))] f(x; \hat{\mu}, \hat{s}^2) dx$$

To summarize, we have

$$d_{TV}(g(x; k, \mu, s^2), f(x; \hat{\mu}, \hat{s}^2)) \leq \max \left(\int_{\hat{\mu}}^{+\infty} [\exp(V(x; k, \hat{\mu})) - 1] f(x; \hat{\mu}, \hat{s}^2) dx, \int_{-\infty}^{\hat{\mu}} [1 - \exp(V(x; k, \hat{\mu}))] f(x; \hat{\mu}, \hat{s}^2) dx \right)$$

As we have shown in Equation (12) of the main manuscript, $\hat{s}^2 = \left(\frac{2k}{\hat{\mu}^2} + C\right)^{-1}$, where $C > 0$. This suggests that $\hat{s} \leq \hat{\mu}/\sqrt{2k}$. Therefore

$$d_{TV}(g(x; k, \mu, s^2), f(x; \hat{\mu}, \hat{s}^2)) \leq \max \left(\int_{\hat{\mu}}^{+\infty} [\exp(V(x; k, \hat{\mu})) - 1] f(x; \hat{\mu}, \hat{\mu}/2k) dx, \int_{-\infty}^{\hat{\mu}} [1 - \exp(V(x; k, \hat{\mu}))] f(x; \hat{\mu}, \hat{\mu}/2k) dx \right)$$

Since $V(x; \mu, s^2)$ and $f(x; \mu, s^2)$ are location-scale families, the above expression can be made free of $\hat{\mu}$ and thus μ and s^2 :

$$d_{TV}(g(x; k, \mu, s^2), f(x; \hat{\mu}, \hat{s}^2)) \leq \max \left(\int_1^{+\infty} [\exp(V(x; k, 1)) - 1] f(x; 1, 1/2k) dx, \int_{-\infty}^1 [1 - \exp(V(x; k, 1))] f(x; 1, 1/2k) dx \right) \quad (\text{A.2})$$

This upper bound on the total variation distance decreases as k increases and it goes to 0 as $k \rightarrow \infty$. This suggests the convergence of the approximating normal distribution to the density family g in total variation sense. We also plot this upper bound as a function of k to verify the conclusion. It is shown in the Figure A.1.

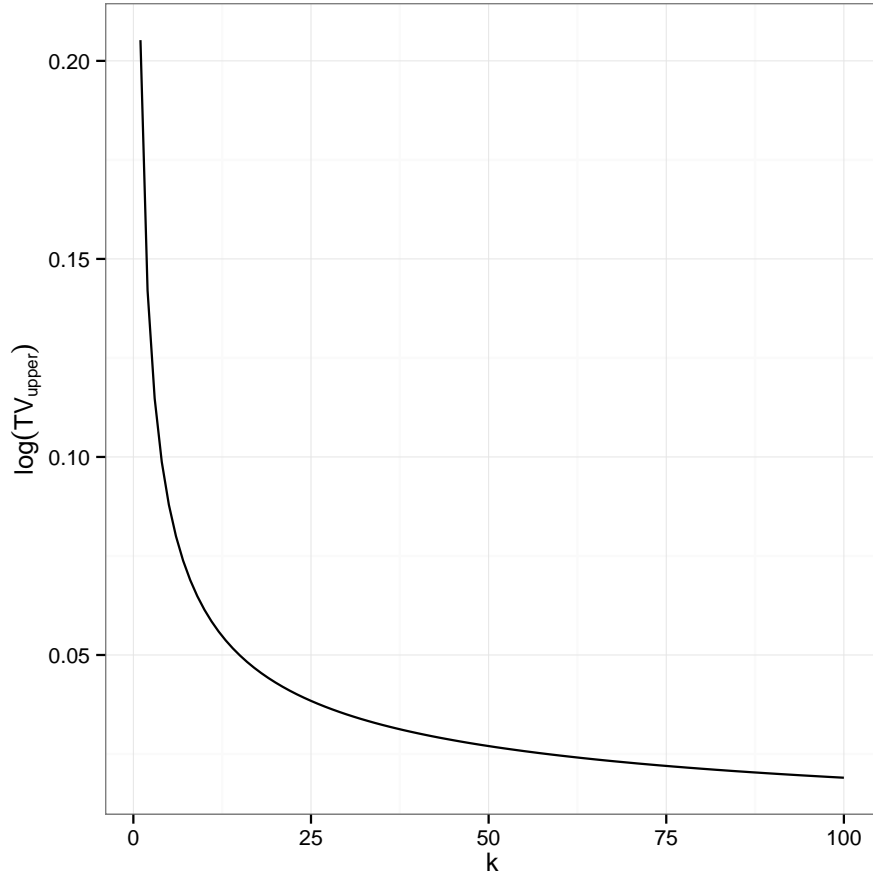


Figure A.1: Upper bound of the total variation distance of Laplace approximation in (11) to the density in (10) as given in (A.2) when frequency k increases.

E Details of self-consistent estimates in section 3.1

First we estimate σ and then we transform the data $n_{i,j}$ into $\sqrt{n_{i,j}/\sigma_i}$. If $n_{i,j}$ is representative and σ is estimated accurately, we have $\sqrt{n_{i,j}/\sigma_i} = c_j Q_{i,j}^+$. If the covariance matrix of \mathbf{Q}_i is Σ , then the covariance matrix of $(\sqrt{n_{i,j}/\sigma_i}, j = 1, \dots, J)$ will be $\tilde{\Sigma} = \Lambda \Sigma \Lambda$ where $\Lambda = \text{diag}\{c_1, \dots, c_J\}$.

It is obvious that $(\sqrt{n_{i,j}/\sigma_i}, j = 1, \dots, J)$ is MVN and the correlation matrix will be the same as the induced correlation matrix from Σ . Methods on identifying the covariance matrix using this truncated dataset are abundant and well-studied. One way to do it is the EM algorithm. This estimated covariance matrix will by no means to be the same as Σ , but the induced correlation matrix will be very close to the true correlation matrix induced by Σ . Hence if our interest is on estimating correlation matrix, we can just treat $(\sqrt{n_{i,j}/\sigma_i}, j = 1, \dots, J)$ as the truncated version of the true \mathbf{Q}_i and proceed.

The EM algorithm should then be derived for the following settings. Let $\mathbf{Q}_i \stackrel{iid}{\sim} MVN(\mathbf{0}, \Sigma)$. Instead of observing I independent \mathbf{Q}_i , we only observe the positive entries in each \mathbf{Q}_i and know the rest of the entries are negative. Denote the observed data vector as $\tilde{\mathbf{Q}}_i$. We want to estimate Σ from the data $\tilde{\mathbf{Q}}_i, i = 1, \dots, I$. A standard EM algorithm can be easily formulated as following:

E-step Get the conditional expectation of full data log likelihood, given the observed data. Define two index sets, $\mathcal{A}_i = \{j | \tilde{Q}_{i,j} > 0\}$ and $\mathcal{B}_i = \{j | \tilde{Q}_{i,j} = 0\}$ and by $Q_{\mathcal{A}_i}$ we meant $(Q_{i,j} | j \in \mathcal{A}_i)$. Denote $\mathcal{A} = \{(i, j) | j \in \mathcal{A}_i, i = 1, \dots, I\}$ and $\mathcal{B} = \{(i, j) | j \in \mathcal{B}_i, i = 1, \dots, I\}$. The E-step function at $t + 1$ iteration is,

$$L(\Sigma | \Sigma_t) = \mathbb{E} \left[-\frac{I}{2} \log |\Sigma| - \frac{1}{2} \text{Tr}(\Sigma^{-1} \sum_i \mathbf{Q}_i \mathbf{Q}_i') | \Sigma_t, Q_{\mathcal{A}} = \tilde{Q}_{\mathcal{A}}, Q_{\mathcal{B}} < 0 \right].$$

Notice this expectation is not easy to calculate in general. We use instead Monte Carlo method to approximate it. We sample K copies of \mathbf{Q}_i from the conditional distribution $(\mathbf{Q}_i | Q_{\mathcal{A}_i} = \tilde{Q}_{\mathcal{A}_i}, Q_{\mathcal{B}_i} < 0)$ where $\mathbf{Q}_i \sim MVN(\mathbf{0}, \Sigma_t)$. This forms a truncated multivariate normal distribution. Sampling from this distribution can be computationally expensive. We use the implementation as in Wilhelm (2015). If we denote by $\mathbf{Q}_i^1, \dots, \mathbf{Q}_i^K$ the K samples of \mathbf{Q}_i , L can be approximated as

$$\hat{L}(\Sigma | \Sigma_t) = -\frac{1}{K} \sum_{k=1}^K \left[\text{Tr}(\Sigma^{-1} \sum_i \mathbf{Q}_i^k (\mathbf{Q}_i^k)') \right] - \frac{I}{2} \log |\Sigma|.$$

M-step We seek to maximize \hat{L} with respect to Σ . Due to a well-known fact on the MLE of covariance matrix of multivariate normal, it is straightforward to get

$$\Sigma_{t+1} = \frac{1}{IK} \sum_{i,k} \mathbf{Q}_i^k (\mathbf{Q}_i^k)'$$

When applying this algorithm to our data, we estimated $\tilde{\mathbf{Q}}_i = (\sqrt{n_{i,j}/\sigma_i}, j = 1, \dots, J)$. A summary of the RV-coefficients between the results from the above algorithm and the truth is shown in Figure A.2. We also compared the results here with those from MCMC simulation in Figure A.2. The estimates of \mathbf{S} from MCMC simulation are always better than those given by this fast algorithm but both perform very well.

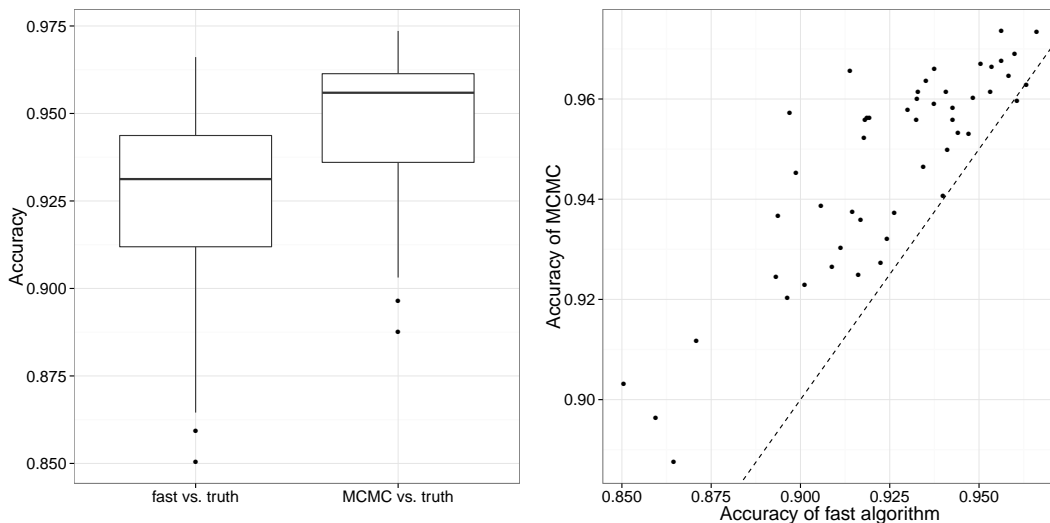


Figure A.2: Left panel: boxplots compare the distribution of RV-coefficients between estimates and truth for our self-consistent algorithm and MCMC simulation. Right panel: scatter plot to indicate per simulation comparison of two algorithms. Dashed line indicates where the two algorithms have identical accuracy.

References

- Abdi, H., A. J. O’Toole, D. Valentin, and B. Edelman (2005). Distatis: The analysis of multiple distance matrices. In *Computer Vision and Pattern Recognition-Workshops, 2005. CVPR Workshops. IEEE Computer Society Conference on*, pp. 42–42. IEEE.
- Anderson, M. J., K. E. Ellingsen, and B. H. McArdle (2006). Multivariate dispersion as a measure of beta diversity. *Ecology Letters* 9(6), 683–693.
- Ando, T. (2009). Bayesian factor analysis with fat-tailed factors and its exact marginal likelihood. *Journal of Multivariate Analysis* 100(8), 1717–1726.
- Barrientos, A. F., A. Jara, F. A. Quintana, et al. (2012). On the support of maceacherns dependent dirichlet processes and extensions. *Bayesian Analysis* 7(2), 277–310.
- Bhattacharya, A. and D. B. Dunson (2011). Sparse Bayesian infinite factor models. *Biometrika* 98(2), 291.

- Brix, A. (1999). Generalized gamma measures and shot-noise cox processes. *Advances in Applied Probability*, 929–953.
- Caporaso, J. G., J. Kuczynski, J. Stombaugh, K. Bittinger, F. D. Bushman, E. K. Costello, N. Fierer, A. G. Pena, J. K. Goodrich, J. I. Gordon, and R. Knight (2010). Qiime allows analysis of high-throughput community sequencing data. *Nature methods* 7(5), 335–336.
- Caporaso, J. G., C. L. Lauber, W. A. Walters, D. Berg-Lyons, C. A. Lozupone, P. J. Turnbaugh, N. Fierer, and R. Knight (2011). Global patterns of 16s rRNA diversity at a depth of millions of sequences per sample. *Proceedings of the National Academy of Sciences* 108(Supplement 1), 4516–4522.
- Carvalho, C. M., J. Chang, J. E. Lucas, J. R. Nevins, Q. Wang, and M. West (2008). High-dimensional sparse factor modeling: applications in gene expression genomics. *Journal of the American Statistical Association* 103(484).
- Daley, D. J. and D. Vere-Jones (1988). An introduction to the theory of point processes.
- DeSantis, T. Z., P. Hugenholtz, N. Larsen, M. Rojas, E. L. Brodie, K. Keller, T. Huber, D. Dalevi, P. Hu, and G. L. Andersen (2006). Greengenes, a chimera-checked 16s rRNA gene database and workbench compatible with arb. *Applied and environmental microbiology* 72(7), 5069–5072.
- Dethlefsen, L., M. McFall-Ngai, and D. A. Relman (2007). An ecological and evolutionary perspective on human–microbe mutualism and disease. *Nature* 449(7164), 811–818.
- Dethlefsen, L. and D. A. Relman (2011). Incomplete recovery and individualized responses of the human distal gut microbiota to repeated antibiotic perturbation. *Proceedings of the National Academy of Sciences* 108(Supplement 1), 4554–4561.
- DiGiulio, D., B. J. Callahan, P. J. McMurdie, E. K. Costello, D. J. Lyell, A. Robaczewska, C. L. Sun, D. S. A. Goltsman, R. J. Wong, G. Shaw, D. K. Stevenson, S. Holmes, and R. D. A. R. (2015). Temporal and spatial variation of the human microbiota during pregnancy. to appear.
- Eren, A. M., G. G. Borisy, S. M. Huse, and J. L. M. Welch (2014). Oligotyping analysis of the human oral microbiome. *Proceedings of the National Academy of Sciences* 111(28), E2875–E2884.
- Escoufier, Y. (1973). Le traitement des variables vectorielles. *Biometrics*, 751–760.
- Ferguson, T. S. (1973). A Bayesian analysis of some nonparametric problems. *The annals of statistics*, 209–230.
- Grice, E. A. and J. A. Segre (2011). The skin microbiome. *Nature Reviews Microbiology* 9(4), 244–253.

- Griffin, J. E., M. Kolossiatis, and M. F. J. Steel (2013). Comparing distributions by using dependent normalized random-measure mixtures. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 75(3), 499–529.
- Hastie, T., R. Tibshirani, B. Narasimhan, and G. Chu (2003). Pam: prediction analysis for microarrays.
- Holmes, I., K. Harris, and C. Quince (2012). Dirichlet multinomial mixtures: generative models for microbial metagenomics. *PloS one* 7(2), e30126.
- Holmes, S. (2008). Multivariate data analysis: the french way. In *Probability and statistics: Essays in honor of David A. Freedman*, pp. 219–233. Institute of Mathematical Statistics.
- Ishwaran, H. and L. F. James (2001). Gibbs sampling methods for stick-breaking priors. *Journal of the American Statistical Association* 96(453).
- James, L. F. (2002). Poisson process partition calculus with applications to exchangeable models and bayesian nonparametrics. *arXiv preprint math/0205093*.
- James, L. F., A. Lijoi, and I. Prünster (2009). Posterior analysis for normalized random measures with independent increments. *Scandinavian Journal of Statistics* 36(1), 76–97.
- Kingman, J. (1967). Completely random measures. *Pacific Journal of Mathematics* 21(1), 59–78.
- Kostic, A. D., D. Gevers, H. Siljander, T. Vatanen, T. Hyötyläinen, A. Hämäläinen, A. Peet, V. Tillmann, P. Pöhö, and I. Mattila (2015). The dynamics of the human infant gut microbiome in development and in progression toward type 1 diabetes. *Cell host & microbe* 17(2), 260–273.
- Lavit, C., Y. Escoufier, R. Sabatier, and P. Traissac (1994). The ACT (statis method). *Computational Statistics & Data Analysis* 18(1), 97–119.
- Lee, S. and X. Song (2002). Bayesian selection on the number of factors in a factor analysis model. *Behaviormetrika* 29(1), 23–39.
- Lijoi, A., R. H. Mena, and I. Prünster (2005). Hierarchical mixture modeling with normalized inverse-gaussian priors. *Journal of the American Statistical Association* 100(472), 1278–1291.
- Lijoi, A., R. H. Mena, and I. Prünster (2007). Controlling the reinforcement in Bayesian non-parametric mixture models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 69(4), 715–740.
- Lijoi, A. and I. Prünster (2010). Models beyond the dirichlet process. In N. L. Hjort, C. Holmes, P. Müller, and S. G. Walker (Eds.), *Bayesian nonparametrics*, Chapter 3, pp. 80–136. Cambridge University Press.

- Lopes, H. F. and M. West (2004). Bayesian model assessment in factor analysis. *Statistica Sinica* 14(1), 41–68.
- Lozupone, C. and R. Knight (2005). Unifrac: a new phylogenetic method for comparing microbial communities. *Applied and environmental microbiology* 71(12), 8228–8235.
- Lucas, J., C. Carvalho, Q. Wang, A. Bild, J. R. Nevins, and M. West (2006). Sparse statistical modelling in gene expression genomics. *Bayesian Inference for Gene Expression and Proteomics* 1.
- MacEachern, S. N. (2000). Dependent dirichlet processes. *Unpublished manuscript, Department of Statistics, The Ohio State University*.
- McMurdie, P. J. and S. Holmes (2013). phyloseq: an r package for reproducible interactive analysis and graphics of microbiome census data. *PLOS one*.
- McMurdie, P. J. and S. Holmes (2014). Waste not, want not: why rarefying microbiome data is inadmissible. *PLoS Comput Biol* 10(4), e1003531.
- Mercer, J. (1909). Functions of positive and negative type, and their connection with the theory of integral equations. *Philosophical transactions of the royal society of London. Series A, containing papers of a mathematical or physical character*, 415–446.
- Muliere, P. and L. Tardella (1998). Approximating distributions of random functionals of ferguson-dirichlet priors. *Canadian Journal of Statistics* 26(2), 283–297.
- Müller, P., F. Quintana, and G. Rosner (2004). A method for combining inference across related nonparametric Bayesian models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 66(3), 735–749.
- Oksanen, J., F. G. Blanchet, R. Kindt, P. Legendre, P. R. Minchin, R. B. O’Hara, G. L. Simpson, P. Solymos, M. H. H. Stevens, and H. Wagner (2015, November). vegan: Community Ecology Package.
- Paulson, J. N., O. C. Stine, H. C. Bravo, and M. Pop (2013). Differential abundance analysis for microbial marker-gene surveys. *Nature methods* 10(12), 1200–1202.
- Peiffer, J. A., A. Spor, O. Koren, Z. Jin, S. G. Tringe, J. L. Dangl, E. S. Buckler, and R. Ley (2013). Diversity and heritability of the maize rhizosphere microbiome under field conditions. *Proceedings of the National Academy of Sciences* 110(16), 6548–6553.
- Press, S. J. and K. Shigemasu (1989). Bayesian inference in factor analysis. In *Contributions to probability and statistics*, pp. 271–287. Springer.

- Quince, C., E. E. Lundin, A. N. Andreasson, D. Greco, J. Rafter, N. J. Talley, L. Agreus, A. F. Andersson, L. Engstrand, and M. D’Amato (2013). The impact of crohn’s disease genes on healthy human gut microbiota: a pilot study. *Gut*, gutjnl–2012.
- Ravel, J., P. Gajer, Z. Abdo, G. M. Schneider, S. S. K. K., S. L. McCulle, S. Karlebach, R. Gorle, J. Russell, C. O. Tacket, and R. M. Brotman (2011). Vaginal microbiome of reproductive-age women. *Proceedings of the National Academy of Sciences* 108(Supplement 1), 4680–4687.
- Regazzini, E., A. Lijoi, and I. Prünster (2003). Distributional results for means of normalized random measures with independent increments. *Annals of Statistics*, 560–585.
- Robert, P. and Y. Escoufier (1976). A unifying tool for linear multivariate statistical methods: the RV-coefficient. *Applied statistics*, 257–265.
- Rodríguez, A., D. B. Dunson, and A. E. Gelfand (2009). Bayesian nonparametric functional data analysis through density estimation. *Biometrika* 96(1), 149–162.
- Rosen, M. J., B. J. Callahan, D. S. Fisher, and S. Holmes (2012). Denoising pcr-amplified metagenome data. *BMC bioinformatics* 13(1), 283.
- Rowe, D. B. (2002). *Multivariate Bayesian statistics: models for source separation and signal unmixing*. CRC Press.
- Turnbaugh, P. J., M. Hamady, T. Yatsunenko, B. L. Cantarel, A. Duncan, R. E. Ley, M. L. Sogin, W. J. Jones, B. A. Roe, J. P. Affourtit, M. Egholm, B. Henrissat, A. C. Heath, R. Knight, and J. I. Gordon (2009, Jan). A core gut microbiome in obese and lean twins. *Nature* 457(7228), 480–484.
- Wilhelm, G. S. with contributions from Manjunath, B. (2015, August). tmvtnorm: Truncated Multivariate Normal and Student t Distribution.